

*Article*

## **Examiner Quality and Consistency across LanguageCert Writing Tests**

**Yiannis Papargyris\***

LanguageCert, Greece

**Zi Yan**

The Education University of Hong Kong, China

Received: 1 March 2020/Accepted: 1 May 2020/Published: 1 July 2020

### **Abstract**

This paper reports on a study of the training and standardisation of examiners who mark LanguageCert's International ESOL (IESOL) suite of English language tests linked to the Common European Framework of Reference (CEFR). Subjects in the study were a set of examiners (N=27) who had been marking LanguageCert's IESOL Writing tests across the six CEFR levels. The focus of the study was on the consistency of marking in terms of severity within and across the six tests that the examiners mark.

Correlations between examiner person measures across all six tests indicated that examiners were broadly consistent across tests, with examiner person measures generally correlating highly with their 'partner' test: A1 with A2, C1 with C2, and B1 with B2 tests. LanguageCert examiners – who undergo careful training and standardisation – may therefore be seen to mark consistently and accurately across a range of ability levels.

### **Keywords**

Examiner quality, examiner consistency, marking, testing, training, LanguageCert

## **1 Introduction**

One of the maxims of assessment is that tests be valid and provide accurate assessments of candidates' abilities: in particular in the context of how far a given test score may be interpreted as an indicator of the abilities or constructs to be measured (Bachman & Palmer, 1996; Messick, 1989). Under such a precondition, the marking of candidates' writing therefore needs to be accurate if reliable assessments are to emerge. However, such accurate marking in performance assessment involving examiner judgment is an enduring challenge because scores assigned to candidate performance are mediated, interpreted, and applied by examiners who are a potential source of error (Engelhard, 2002). From this, it naturally

---

\*Corresponding Author. Email: [yiannis.papargyris@peoplecert.org](mailto:yiannis.papargyris@peoplecert.org)

follows that examiners need to be properly trained and standardised – in particular with subjectively-marked performance tests such as Speaking and Writing.

This paper reports on a study of the training and standardisation of examiners who mark *LanguageCert's* International ESOL (IESOL) suite of English language tests linked to the Common European Framework of Reference (CEFR). Subjects in the study were a set of examiners (N=27) who had been marking *LanguageCert's* IESOL Writing tests across the six CEFR levels. The focus of the study was on the consistency of marking in terms of severity within and across the six tests that the examiners mark.

## 2 Background to the Tests, Examiners, Scripts

The data in the study were drawn from six examinations which comprise *LanguageCert's* International ESOL suite of English language tests. In the LanguageCert Writing tests, candidates complete two writing tasks which elicit a range of writing skills. Responses are marked using an analytic mark scheme which reflects the CEFR descriptors. Separate marks are awarded by marking examiners for different aspects of writing ability – Task Fulfilment, Accuracy and Range of Grammar, Accuracy and Range of Vocabulary, and Organisation of the text. This set of criteria ensures that a wide range of writing skills are considered, thus enhancing the reliability and representativeness of test scores.

The format of the tests and the nature of the assessment criteria reflect the broad multi-faceted construct underlying these examinations. Communicative ability is the primary concern, while accuracy and range become increasingly important as the CEFR level of the test increases.

### 2.1 Examiner Training

The importance of examiner training in any English language examination is an issue which has long been accepted as an essential factor in determining the reliability of a test (see e.g., [Webb et al., 1990](#)). Although empirical studies on examiner training have generated mixed results, a general consensus is that examiner training, if well designed, can improve the reliability and validity of examiner-mediated assessment ([Kang et al., 2019](#)). Studies have shown trained examiners to be more reliable ([Saito, 2008](#)) as well as more self-consistent ([Davis, 2016](#)) than untrained examiners.

In the case of performance-based assessment, it is important to attempt to ensure reliability through extensive examiner training and standardisation, including even sanctioning inconsistent examiners (see [Elder et al., 2007](#)).

[Webb et al. \(1990\)](#) discuss the problems associated with examiner stringency, leniency and inconsistency. They state that problems with examiner stringency and leniency can be handled by statistical adjustment. They make it clear nonetheless that examiner training is essential for other problems – specifically, examiner inconsistency. As [Weigle \(1998\)](#) notes, examiner training was more effective in enhancing intra-examiner reliability than inter-examiner reliability. [Lumley and McNamara \(1995\)](#), in discussing inconsistency in examiners, report that training and standardisation are not only essential, but also that further moderation is required shortly before the administration of Writing or Speaking Tests because a time gap between the training and the assessment event reveals that inconsistencies re-emerge.

In order to address the issue of consistency, severity, and leniency amongst the group of LanguageCert examiners, Multi-Faceted Rasch Analysis (MFRA), via the computer program FACETS ([Linacre, 2020](#)) has been utilised. A brief outline of the Rasch measurement model and MFRA is given below.

## 2.2 The Rasch Model

The use of the Rasch model enables different facets (person ability, examiner severity, and item difficulty in the current instance) to be modelled together. First, in the standard Rasch model, the aim is to obtain a unified and interval metric for measurement. The Rasch model converts ordinal raw data into interval measures which have a constant interval meaning and provide objective and linear measurement from ordered category responses (Wright, 1997). This is not unlike measuring length using a ruler, with the units of measurement in Rasch analysis (referred as the ‘logit’) evenly spaced along the ruler. Second, once a common metric is established for measuring different phenomena (candidates and test items being the most obvious), person ability estimates are independent of the items used, with item difficulty estimates being independent of the sample recruited because the estimates are calibrated against a common metric rather than against a single test situation (for person ability estimates) or a particular sample of candidates (for item difficulty estimates). Third, Rasch analysis prevails over Classical Test Analysis statistics by calibrating persons and items onto a single unidimensional latent trait scale (Bond et al., 2020).

Person measures and item difficulties are placed on an ordered trait continuum by which direct comparisons between person measures and item difficulties can be easily conducted. Consequently, results can be interpreted with a more general meaning. The use of MFRA adds flexibility to the measurement by allowing the incorporation of facets in addition to person ability and item difficulty. As the current study focuses on the examiner severity facet (leniency vs stringency of marking) in IESOL Writing tests, the MFRA analysis includes three facets: candidates, items, and examiners.

## 2.3 Principles and Procedures in Training Examiners

As stated earlier, in any examination of direct performance it is important to attend to the question of examiner reliability. Although there is no agreement regarding the most effective training and standardisation methods (Kogan et al., 2015), in assessments of performance which rely wholly on examiner applications of the criteria established for the assessment, reliability can be established through a process of:

- agreement on the validity of assessment constructs
- creation of detailed specifications
- creation of valid, detailed, and usable descriptors
- provision of credible and regular examiner training
- standardisation

See also Feldman et al. (2012), where a cogent summary of different modes of examiner training is provided.

The purpose of standardising examiners is to ensure that strong measures of agreement occur whenever a number of examiners apply grade descriptors to a criterion-referenced assessment instrument. This is the case with the *LanguageCert* Writing tests. In criterion-referenced assessment, which depends on the application of examiners’ judgements to the criteria described in the descriptors, it is important that two principles are adhered to:

- Judgements by one examiner over time with a number of candidates need to be consistent.
- Different examiners judging an individual candidate should provide assessments that are in close agreement.

There are a number of well-established standard procedures that can be used to train and standardise language examiners (see e.g., Coniam & Falvey, 2018). These procedures were applied in the specific training procedures used with examiners for the IESOL Writing Tests and are described below.

## 2.4 Participants

All writing examiners must meet minimum requirements in terms of professional qualifications and experience in order to be eligible for consideration as an examiner. Prospective examiners go through a training process before they are approved and allowed to mark. The training process includes marking sample scripts. Candidates for the examiner role must show they can mark accurately and consistently before they are certificated as examiners. During live marking, where an examiner is found to be marking inaccurately and/or inconsistently, they may be removed from the marking session and/or retrained, or dismissed as an examiner. Examiners are then monitored on an ongoing basis and required to attend standardisation meetings on a regular basis.

Participants involved 27 examiners who have been marking *LanguageCert*'s IESOL suite of examinations for a considerable period of time. All 27 examiners marked the A1 and A2 scripts, whereas only 24 examiners were eligible at the time to conduct assessments at the other four CEFR levels, i.e., B1 to C2, as a result of the *LanguageCert* examiner training process.

## 2.5 Standardisation

Examiners were familiar with the rating scales since they had been using them for five years. The standardisation session described in this paper took place in 2018 and is a regular feature of re-training and standardising undergone by *LanguageCert* assessment personnel. The process was led by the Chief Examiner, who has marked examinations linked to the CEFR for over 20 years.

Examiners were first given the rating scales, and *LanguageCert*'s *Guide for Examiners*, and asked to familiarise themselves with the constructs and levels in the scales. Some brief discussion was then followed by two stages of training, Induction and Training, each consisting of the assessment of 36 benchmarked scripts – six per CEFR level – and subsequent discussion of: queries; potential discrepancies between raters; the applicability of descriptors, etc. The sample scripts shared with examiners during the Induction and Training stages exemplified the four criteria along with the performance descriptors which constitute the marking scheme.

Over a period of a day and a half, examiners then marked, one test at a time, six scripts from each of the six tests in the *LanguageCert* IESOL suite (i.e., from A1 to C2). The marking began with the six A1 tests, progressing upwards. After each set of marking and after all examiners had submitted their awarded marks, the Chief Examiner revealed the scores he had awarded and led a discussion of the merits of different scripts.

*LanguageCert* training and standardisation procedures and practices may be seen therefore to equate with those employed primarily under a performance dimension training (PDT) (see [Kogan et al., 2015](#)) as all three training stages (Induction, Training, Standardisation) are based on the assessment of a series of sample scripts (performances), selected and/or adapted to demonstrate certain issues in candidate performance. To account for potential discrepancies in marking as a result of raters' idiosyncratic tendencies (e.g., excessive leniency/severity), elements of a frame of reference training (FoRT) methodology were employed so that the role of subjectivity in the application of the marking criteria was minimised.

## 2.6 The IESOL Writing Test

The IESOL Writing tests comprise two tasks, as laid out in Table 1.

Table 1  
*IESOL Writing Test Tasks and Scales*

Level	Part 1 : Candidates produce	Word length	Part 2 : Candidates produce	Word length
A1	four sentences on a specified topic	30	a simple text for a specified reader	20-30
A2	an informal response to an informal text	30-50	a neutral response to a specified public reader	30-50
B1	a neutral or formal text for a public audience	70-100	a letter using informal language	100-120
B2	a neutral or formal text for a public audience	100-150	a text using informal language	150-200
C1	a neutral or formal text for a public audience	150-200	a text using informal language	250-300
C2	a neutral or formal text for a public audience	200-250	a text using informal language	250-300

Concerning marking, all tasks conform to CEFR ‘can do’ statements for writing and are assessed on a four-point scale on four domains, as Table 2 illustrates.

Table 2  
*Rating Scale Domains*

Task Fulfilment
Accuracy and range of grammar
Accuracy and range of vocabulary
Organisation

### 3 Method

The key research question for this study is whether examiner severity will be comparable within each test and across tests at the six CEFR levels, i.e., whether examiners will apply the marking descriptors accurately and be consistently lenient / severe on tests within a level and across levels. Two indicators of examiner severity and consistency were examined to address the research question.

The first indicator generated from the Rasch analysis is the person fit statistic. This statistic is not a direct indicator but a pre-requisite of examiner consistency. Examiner performance has to satisfy Rasch measurement requirements (i.e., the fit to the Rasch model) before any meaningful discussions of severity estimates may be made. The computer program FACETS (Linacre, 2020) provides a number of statistics which give an indication as to how well the data fits the model. One of these is the mean square statistic. For person fit statistics (examiners, in our case), acceptable practical limits of fit have been proposed as 0.5 for the lower limit and 1.5 for the upper limit (Lunz & Stahl, 1990).

The second indicator relates to examiner invariance across tests. While MFRA provides a framework for obtaining fair measurements of examinee ability that can be statistically invariant over examiners, tasks, and other aspects of performance assessment procedures, this only applies across one test. In the current study, examiner invariance across the six tests is examined via the Spearman’s rho, which reports rank order correlations between tests. A high correlation indicates consistency of rank order of examiner severity estimates.

## 4 Results and Discussion

### 4.1 Examiner Fit to the Rasch Model

As the cornerstone of good rating is fit to the Rasch model, results are first presented below for the examiners on each of the six tests. Tables 2a and 2b present the results for the 27 examiners who participated in the standardisation exercise. As mentioned, 24 examiners marked all six tests, with the whole cohort of 27 examiners marking tests A1 and A2. In the tables, infit is reported. Infit shows the ‘big picture’ in that it scrutinises the internal structure of a facet (examiners, in this case). Generally speaking, high infit (above 1.5) values are more critical and suggest an examiner’s ratings to be rather ‘scattered’, providing a confused picture about the placement of the examiner’s ratings. Very small (below 0.5) infit values indicate only very small variation in the data, thereby providing little information to articulate clear and meaningful judgments about the examiner, and their ratings.

In Tables 3 and 4 below, infit figures above 1.5 are highlighted in yellow, and infit figures below 0.5 highlighted in green. In the data and discussion that follows, all examiner names have been anonymised.

Table 3

*Examiner Measures for Tests A1 and A2 (N=27)*

Examiners	Nu	A1-Measure	A1-S.E.	A1-Infit	A2-Measure	A2-S.E.	A2-Infit
Andy	1	-0.1	0.46	0.64	0.69	0.44	1.14
Brian	2	0.5	0.44	0.85	1.47	0.45	0.87
Cathy	3	-0.78	0.5	0.68	-0.92	0.47	1
Dot	4	0.31	0.44	0.52	-0.09	0.45	1.29
Ellen	5	0.69	0.43	0.65	0.3	0.44	0.71
Fred	6	0.31	0.44	0.81	-0.09	0.45	0.72
Gary	7	0.11	0.45	1.7	-1.15	0.48	1.06
Terri	8	-1.61	0.56	0.61	-0.92	0.47	0.86
Iris	9	0.11	0.45	0.93	-0.09	0.45	0.88
Jack	10	1.92	0.41	1.01	0.69	0.44	0.76
Katie	11	0.88	0.43	1.43	0.11	0.45	1
Lenny	12	0.31	0.44	1.11	-0.92	0.47	1.05
Martha	13	-1.61	0.56	0.61	-0.92	0.47	0.86
Nonie	14	-0.54	0.48	0.53	0.11	0.45	0.61
Oliver	15	0.31	0.44	0.94	-1.62	0.5	0.84
Perry	16	0.5	0.44	0.99	1.08	0.44	1.03
Queenie	17	-1.04	0.52	2.46	-0.09	0.45	0.83
Robert	18	0.31	0.44	1.17	0.69	0.44	1.7
Susan	19	-0.54	0.48	1.21	-0.5	0.46	1
Terri	20	-1.31	0.54	0.76	-0.29	0.45	0.83
Ursula	21	-0.1	0.46	0.77	-0.5	0.46	0.78
Vanesa	22	0.11	0.45	1.53	1.08	0.44	1.39
Windy	23	-0.1	0.46	1.27	0.5	0.44	1.54
Xerxes	24	0.31	0.44	1.01	0.69	0.44	0.79
Yana	25	1.24	0.42	0.68	1.28	0.44	0.61
Zoe	26	-0.1	0.46	1.2	-0.71	0.46	0.98
Albert	27	-0.1	0.46	0.9	0.11	0.45	0.96

Table 4  
*Examiner Measures for Tests B1, B2, C1 and C2 (N=24)*

Examiners	Nu	B1-Measure	B1-S.E.	B1-Infit	B2-Measure	B2-S.E.	B2-Infit	C1-Measure	C1-S.E.	C1-Infit	C2-Measure	C2-S.E.	C2-Infit
Andy	1	0.88	0.41	0.67	1.82	0.46	0.81	0.63	0.36	1.42	0.75	0.43	0.88
Brian	2	0.54	0.41	0.96	1.19	0.46	0.68	0.24	0.36	0.64	1.29	0.43	0.84
Cathy	3												
Dot	4	0.54	0.41	1.16	-0.23	0.45	1.11	0.63	0.36	1.3	0	0.44	1.04
Ellen	5	-1.41	0.49	0.81	-0.23	0.45	0.88	0.63	0.36	0.71	0.75	0.43	0.67
Fred	6	-0.96	0.47	1.11	-0.64	0.45	0.81	0.5	0.36	0.91	1.29	0.43	1.27
Gary	7	-1.9	0.51	1.59**	-0.84	0.46	1.04	-0.98	0.38	0.97	-0.38	0.44	1.47
Terri	8												
Iris	9	-1.65	0.5	1.77	-0.43	0.45	0.87	0.11	0.36	1.11	-0.57	0.44	0.79
Jack	10	0.71	0.41	1.07	1.61	0.46	0.94	0.5	0.36	0.42*	0.93	0.43	0.58
Katie	11	-0.54	0.45	0.67	-0.03	0.45	0.88	0.24	0.36	0.57	0.56	0.43	0.46*
Lenny	12	-0.16	0.43	1.23	-0.03	0.45	1.04	-0.98	0.38	1.27	-1.77	0.46	0.44*
Martha	13												
Nonie	14												
Oliver	15	-1.18	0.48	1.44	-1.05	0.46	0.63	-0.7	0.37	0.78	-0.96	0.45	0.56
Perry	16	0.2	0.42	0.97	-0.43	0.45	1.27	-0.98	0.38	1.3	-1.15	0.45	0.72
Queenie	17	-0.75	0.46	0.7	0.18	0.45	0.66	-0.43	0.37	1.05	-1.15	0.45	0.95
Robert	18	0.54	0.41	1.18	0.58	0.45	1.15	0.11	0.36	1.59**	0.93	0.43	1.22
Susan	19	0.2	0.42	0.83	-0.84	0.46	0.93	-0.7	0.37	1.28	-1.36	0.45	1.4
Terri	20	0.71	0.41	0.6	0.38	0.45	0.72	1.53	0.36	1.02	1.47	0.42	0.36*
Ursula	21	0.71	0.41	0.6	-0.43	0.45	2.09**	-0.43	0.37	1.02	0.38	0.43	1.14
Vanesa	22	1.04	0.41	0.84	0.99	0.45	1.02	-0.29	0.37	1.12	0.38	0.43	1.67**
Windy	23	0.02	0.43	1.14	-0.03	0.45	0.77	-0.29	0.37	1	0.56	0.43	2.03**
Xerxes	24	1.85	0.4	0.66	-0.84	0.46	0.33	0.24	0.36	0.59	0.38	0.43	0.54
Yana	25	1.04	0.41	0.77	-0.23	0.45	1.03	0.5	0.36	0.89	-0.57	0.44	0.74
Zoe	26	-1.18	0.48	1.48	-0.64	0.45	1.63**	-0.43	0.37	0.85	-1.77	0.46	0.71
Albert	27	0.71	0.41	0.76	0.18	0.45	1.09	0.37	0.36	0.66	0	0.44	1.48

\*. Correlation significant at the 0.05 level. \*\*. Correlation significant at the 0.01 level

As can be seen from the data in the above table, examiner fit to the model was generally good; there were only one or two examiners who exhibited high infit (i.e., with a mean square of over 1.5) on different tests. From the total of 150 examiner/test infit results, there are 10 instances of infit greater than 1.5 and five instances of infit lower than 0.5. More than one instance of poor fit was observed with only three examiners – Gary, Robert, Windy.

## 4.2 Examiner Consistency Across Tests

Having established that examiners broadly fit the model, the next step involves investigating examiner consistency across tests. Table 5 presents the results of rank order correlations (via Spearman's rho) conducted against examiner person measures across the 6 tests. Correlations significant at the 0.01 level are highlighted in yellow, and those significant at the 0.05 level in green.

Table 5

*Examiner Measure Rank Order Correlations Across the Six Tests*

CEFR level	Correlation Details	A1 Measure	A2 Measure	B1 Measure	B2 Measure	C1 Measure	C2 Measure
A1 Measure	Correlation Coefficient	--	.531**	.014	.035	.183	.187
	Sig. (2-tailed)	.	.004	.951	.876	.404	.393
	N	27	27	23	23	23	23
A2 Measure	Correlation Coefficient	.531**	--	.582**	.551**	.395	.425*
	Sig. (2-tailed)	.004	.	.004	.006	.062	.043
	N	27	27	23	23	23	23
B1 Measure	Correlation Coefficient	.014	.582**	--	.459*	.388	.302
	Sig. (2-tailed)	.951	.004	.	.028	.067	.162
	N	23	23	23	23	23	23
B2 Measure	Correlation Coefficient	.035	.551**	.459*	--	.447*	.514*
	Sig. (2-tailed)	.876	.006	.028	.	.033	.012
	N	23	23	23	23	23	23
C1 Measure	Correlation Coefficient	.183	.395	.388	.447*	--	.696**
	Sig. (2-tailed)	.404	.062	.067	.033	.	.000
	N	23	23	23	23	23	23
C2 Measure	Correlation Coefficient	.187	.425*	.302	.514*	.696**	--
	Sig. (2-tailed)	.393	.043	.162	.012	.000	.
	N	23	23	23	23	23	23

\*. Correlation significant at the 0.05 level. \*\*. Correlation significant at the 0.01 level

As may be seen from Table 5, in general, tests (that is, via examiner person measures) correlate highly with their ‘partner’: hence the A1 and A2 tests correlate highly (at the  $p < .01$  level), as do the C1 and C2 tests; and the B1 and B2 tests correlate quite highly (at the  $p < .05$  level). While the A2 test appears to correlate with almost all tests, all tests correlate quite highly with at least two or more different tests. The implication of these correlations is that the rank order of the examiners is broadly consistent across tests: if an examiner is going to be strict on one test, it is quite likely that they will be strict on other tests.

## 5 Conclusion

This study has examined the issue of examiner severity and invariance across *LanguageCert*’s six CEFR-linked IESOL Writing tests. The research question was whether examiner severity would be comparable across the six tests, i.e., examiners would be consistently severe on each test. If examiners are seen to be erratic in their severity at some levels but not at others, this may impact on fairness in terms of grades awarded to candidates.

An examination of 27 examiners standardised to mark *LanguageCert*’s six CEFR-linked IESOL Writing tests, illustrated that examiner fit to the Rasch model was generally good – a key background consideration.

From correlations run among the examiner person measures across all six tests, a rank order emerged indicating that examiners were broadly consistent across tests. Examiner person measures generally correlated highly with their ‘partner’ test: A1 with A2, C1 with C2, and B1 with B2 tests. While the A2

test correlated with almost all tests, all tests correlated quite highly with at least two or more different tests.

A major implication which arises regarding consistency is the following: if an examiner is going to be strict at one level, they will quite likely be strict at other levels – and strictness can be compensated for by statistical adjustment if necessary. Given that *LanguageCert* examiners undergo careful training and standardisation, what the current study illustrates is that *LanguageCert* examiners may be seen to mark consistently and accurately across a range of ability levels.

## References

- Bachman, L. F., & Palmer, A. S. (1996). *Language testing in practice: Designing and developing useful language tests*. Oxford University Press.
- Bond, T. G., Yan, Z., & Heene, M. (2020). *Applying the Rasch model: Fundamental measurement in the human sciences* (4th ed.). Routledge. <https://doi.org/10.4324/9780429030499>
- Coniam, D., & Falvey, P. (2018). *High-stakes testing: The impact of the LPATE on English language teachers in Hong Kong*. Springer Nature. <https://doi.org/10.1007/978-981-10-6358-9>
- Davis, L. (2016). The influence of training and experience on rater performance in scoring spoken language. *Language Testing*, 33(1), 117-135. <https://doi.org/10.1177/0265532215582282>
- Elder, C., Barkhuizen, G., Knoch, U., & Randow, J. (2007). Evaluating rater responses to an online training program for L2 writing assessment. *Language Testing*, 24(1), 37-64. <https://doi.org/10.1177/0265532207071511>
- Engelhard, G., Jr. (2002). Monitoring raters in performance assessments. In G. Tindal, & T. Haladyna (Eds.), *Large-scale assessment programs for all students: Development, implementation, and analysis* (pp. 261–287). Erlbaum.
- Feldman, M., Lazzara, E. H., Vanderbilt, A. A., & DiazGranados, D. (2012). Rater training to support high-stakes simulation-based assessments. *Journal of Continuing Education in the Health Professions*, 32(4), 279-286. <https://doi.org/10.1002/chp.21156>
- Kang, O., Rubin, D., & Kermad, A. (2019). The effect of training and rater differences on oral proficiency assessment. *Language Testing*, 36(4), 481–504. <https://doi.org/10.1177/0265532219849522>
- Kogan, J. R., Conforti, L. N., Bernabeo, E., Iobst, W., & Holmboe, E. (2015). How faculty members experience workplace-based assessment rater training: A qualitative study. *Medical Education*, 49(7), 692-708. <https://doi.org/10.1111/medu.12733>
- Linacre, J. M. (2020). *Facets computer program for many-facet Rasch measurement*. Winsteps.com. <https://www.winsteps.com/index.htm>
- Lumley, T., & McNamara, T. (1995). Examiner characteristics and examiner bias: Implications for training. *Language Testing*, 12(1), 54-71. <https://doi.org/10.1177/026553229501200104>
- Lunz, M., & Stahl, J. (1990). Judge consistency and severity across grading periods. *Evaluation and the Health Profession*, 13, 425-444. <https://doi.org/10.1177/016327879001300405>
- McNamara, T. (1996). *Measuring second language performance*. Longman.
- Rasch, G. (1960). *Probabilistic models for some intelligence and achievement test*. The University of Chicago Press.
- Messick, S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3 ed., pp. 13-103). American Council on Education.
- Saito, H. (2008). EFL classroom peer assessment: Training effects on rating and commenting. *Language Testing*, 25(4), 553-581. <https://doi.org/10.1177/0265532208094276>

- Webb, L., Raymond, M., & Houston, W. (1990). *Examiner stringency and consistency in performance assessment*. Paper presented at the Annual Meeting of the American Educational Research Association (Boston, MA, April 16-20, 1990).
- Weigle, S. (1998). Using FACETS to model examiner training effects. *Language Testing*, 15(2), 263-287. <https://doi.org/10.1177/026553229801500205>
- Wright, B. D. (1997). A history of social science measurement. *Educational Measurement: Issues and Practice*, 16(4), 33-45. <https://doi.org/10.1111/j.1745-3992.1997.tb00606.x>

**Yiannis Papargyris** is Language Assessment Development Director at LanguageCert. He is an education management professional with over 15 years' experience in the fields of English-medium higher education, qualification development and educational assessment. He is responsible for the development of the LanguageCert exams portfolio.

**Zi Yan** is Associate Professor at the Department of Curriculum and Instruction, the Education University of Hong Kong. His main publications and research interest focus on two related areas: educational assessment and measurement (in particular the application of the Rasch model in educational and psychological research). He is Chief Investigator of several competitive external grants supported by the University Grants Committee (Hong Kong).