

1 **DEVELOPMENTS IN VALIDITY RESEARCH IN**  
2 **SECOND LANGUAGE PERFORMANCE TESTING**

3  
4 *Hacer Hande Uysal\**  
5

6 **ABSTRACT**  
7

8 The present paper aims to provide a short historical overview of the  
9 theoretical developments in validity research in second language performance  
10 testing. A comparative description and critical evaluation of different views  
11 such as the “Trinitarian approach” versus the construct validity model;  
12 “uniform approach,” versus “unified approach” as well as alternative and  
13 critical approaches to validation in L2 performance testing are presented.  
14 These various theoretical approaches are introduced in terms of their  
15 definitions of the validity concept, their suggested requirements for the  
16 validity research, and their attitudes towards reliability and theory while  
17 making interpretations of test scores. The paper also focuses on the current  
18 problems with the applicability of these theoretical approaches, and discusses  
19 future directions in validity research.

20 Key words: Second language assessment, performance assessment,  
21 validity, reliability, validity research  
22

23 **EARLIER THEORIES OF VALIDITY**  
24

25 In earlier times, validity was described as whether the test measures what  
26 it is supposed to measure (Lado, 1961 in Chapelle, 1999). Subject matter  
27 experts used to decide about the quality of the test based on the test content  
28 examining whether test tasks cover a representative sample of the target  
29 domain (Chapelle, 1999). However, this approach was based on a subjective  
30 judgment focusing only on the test content without considering the test scores.  
31 Therefore, while this approach offered evidence to support the domain  
32 relevance and representativeness of the test, it did not provide evidence about  
33 the inferences that could be made from the test scores (Bachman, 1990;  
34 Messick, 1994).

35 Later, Oller (1979, in Chapelle, 1999)<sup>1</sup> put reliability of test scores at the  
36 center of validation. Validity was defined in terms of the degree of the  
37 correlation of test scores with an older or well-established test or criterion  
38 focusing on *criterion related validity*. According to the correlations between  
39 future or present performance and the criterion, criterion validity was later

---

\* Gazi University, Turkey

## DEVELOPMENTS IN VALIDITY RESEARCH IN SECOND LANGUAGE PERFORMANCE TESTING

1 divided into predictive or concurrent validities (Cronbach & Meehl, 1955).  
2 However, this approach was also problematic, because it was not easy to find  
3 a well-defined valid criterion measure all the time; and even if it was found,  
4 validity of this established criterion would also be questionable. Therefore,  
5 criterion based model was not useful in many contexts (Kane, 2001).

6 In 50's, the construct validity was introduced as an alternative to content  
7 and criterion validity, and became one of the several types of validities.  
8 Construct validity was tied to theoretical constructs and started to be  
9 investigated by testing hypotheses related to how well the scores satisfy the  
10 theory (Chapelle, 1999; Kane, 2001). Between 50's and 70's, there were  
11 many kinds of validities –The Trinitarian Model (Shepard, 1993), and while  
12 performing validity research, the type of the validity to be addressed was  
13 chosen according to the purpose of the assessment. (e.g. content validity for  
14 achievement tests; criterion validity for selection and placement decisions,  
15 and the construct validity for theory-based proficiency tests). However, at  
16 those times, validation was still seen as “*one time activity*” (Bachman, 1990).

17 With the development of the construct validity model, limitations of other  
18 validation efforts started to be more apparent (Kane, 2001). The Trinitarian  
19 model was criticized as being “fragmented and incomplete” excluding “score  
20 meaning and social values from test interpretation and use” (Messick, 1995, p.  
21 741). Cronbach & Meehl (1955), for the first time, regarded validity as a  
22 unitary concept including content, criterion, and construct based evidence  
23 under the name of “construct validity.” During 80's, the Trinitarian validity  
24 definition was replaced with a single unified view of validity in the testing  
25 standards (APA, 1985). The focus of interest changed from validating test or  
26 test scores to validating proposed *interpretation of the scores* (Kane, 2001). In  
27 addition, validation became an *on-going process* through which a variety of  
28 empirical evidence about test interpretation and use had to be collected  
29 (Bachman, 1990). In addition, the consequential aspects of validity including  
30 washback, ethics and social responsibility were introduced into validity  
31 discussions (Messick, 1989). Messick (1989) proposed a “progressive  
32 matrix,” which suggested that to *justify* a test score; evidence for construct  
33 validity should be gathered with consideration of value implications of the  
34 interpretation. To *use* the test scores; however, the relevance of the particular  
35 use and social consequences should also be considered.

36 The uniform construct validity approach suggested that interpretations of  
37 all tests – including performance tests – should be validated in the same way  
38 in terms of the theoretical constructs. Messick (1994) stated that adjusting  
39 validity criteria for language performance assessments might cause de-  
40 emphasis on important validity aspects such as construct representativeness  
41 and relevance. Bachman (2002 a,b), although he distinguished between a  
42 construct-based and a task-based approach, suggested that a cognitively based

1 model of language ability and use should be established for all types of  
2 assessment; only then, the similarity between the TLU domain language use  
3 tasks and assessment tasks; adequateness of the domain sampling; and  
4 extrapolation would gain meaning in validation.

5 According to Kane (2001) however, insistence on the necessity of a theory  
6 base for all types of assessments was meaningless especially in the areas  
7 where there is little theory, and the uniform approach to validation that is too  
8 theoretical and ambiguous caused confusion understanding what construct  
9 validity and validation study meant. Although Cronbach (1988) had made an  
10 attempt by suggesting a strong program (necessitated theory, but inapplicable)  
11 and a weak program (abstract, practical, and allowing the use of all kinds of  
12 relevant evidence without any criteria) for validation to reduce this ambiguity,  
13 the “strong” and “weak” arguments were still found to be far from being  
14 definitive and adequate to support the constructs (McNamara, 1996).

## 15 16 **ALTERNATIVE AND CRITICAL APPROACHES: REJECTION** 17 **OF THE THEORY** 18

19 While the problems with regard to the inapplicability of the strong  
20 approach and the lack of criteria in the weak approach continued to cause  
21 ambiguities in validation attempts, new perspectives such as alternative  
22 paradigm and critical theory were included in the validity discussions. If we  
23 look at the validity discussions as a continuum, the “construct validity as a  
24 uniform approach” that requires a strong theory for all assessment types  
25 represents one end and the alternative approaches to validation that are  
26 skeptical of any kind of theory explaining human performance represent the  
27 other. This alternative view demanded for different criteria for validity  
28 judgments in performance assessments claiming that the complex constructs  
29 in human performance cannot be captured by any traditional theories (Lynch,  
30 2001; Moss, 1994).

31 For example, Moss (1994) suggested an alternative hermeneutics  
32 approach to reliability and validity of interpretations. This view argued that  
33 validity was possible without reliability, and did not see inconsistencies in  
34 performances across tasks and among raters as a problem. According to Moss,  
35 it was possible to make generalizations across tasks by developing holistic,  
36 integrative and coherent interpretations based on a *collection of performances*.  
37 Generalization across raters, on the other hand, could be achieved through a  
38 critical dialogue and debate among raters in which initial disagreements  
39 would be resolved, and more refined interpretations would be formed by  
40 considering multiple perspectives and justifying the decisions.

41 Another alternative view was critical language testing, which put  
42 consequential validity at the center of the validity argument. It was suggested

## DEVELOPMENTS IN VALIDITY RESEARCH IN SECOND LANGUAGE PERFORMANCE TESTING

1 that constructs are indefinable as there are multiple perspectives and no truth.  
2 Besides, all tests are subjective, relative, dependent on context, and power  
3 related; thus, there is no true score to be estimated (Shohamy, 2001; Lyncey,  
4 2001). Validity framework according to this view was based merely on  
5 consequences; therefore, information about fairness, ontological authenticity,  
6 cross-referential authenticity, consequential validity, and evolved power  
7 relationships had to be collected (Lynch, 2001).

### 8 9 **A MIDDLE WAY: OBSERVABLE TRAITS VS. THEORETICAL** 10 **CONSTRUCTS**

11  
12 The uniform construct validity approach was too vague and inapplicable,  
13 whereas the alternative-critical approaches focused too much on consequences  
14 and completely rejected constructs and other important validity requirements.  
15 Kane (2001), on the other hand, suggested an alternative approach that was  
16 *unified*, but flexible in which, different kinds of validity arguments to support  
17 different kinds of inferences could be made according to the context. While  
18 the details of the validity argument for each interpretive argument would be  
19 unique, the general approach to specify and evaluate the inferences would be  
20 consistent or unified. Kane's validity definition did not require a theory; yet, it  
21 was still reflecting on the general principles in the construct model. Kane  
22 suggested adopting an argument-based approach – an *interpretative argument*  
23 – rather than validation research.

24 In the interpretative argument, there were several chains of inferences to  
25 be followed: 1) evaluation of performance on each task and giving a score; 2)  
26 generalization of the score beyond the observed to a universe of possible  
27 performances on similar tasks under similar circumstances; 3) extrapolation of  
28 the results beyond the testing context to various other contexts and task  
29 formats; 4) explanation and decision-making based on the theory. Kane  
30 suggested that all evidence relevant to each inference should be collected,  
31 alternative interpretations should be eliminated, and the most problematic  
32 assumptions – the weakest link—should be evaluated (Kane, et. al, 1999;  
33 Kane, 2001). Therefore, generalizability link in performance assessments  
34 should be handled carefully because it is the weakest link due to the use of  
35 small number of tasks representing a narrow range of TLU domain and due to  
36 the variability associated to raters, task, and especially person-task interaction  
37 (Mc Namara, 1997). According to Kane, if generalizability link fails, it is not  
38 possible to talk about extrapolation, and failure of any of the inferences fails  
39 the argument as a whole (Kane, 2001).

40 Kane (2001) also suggested that a distinction should be made between  
41 *theoretical constructs* and *observable attributes*. According to Kane,  
42 theoretical constructs and observable attributes are different both in terms of

1 validity definitions and interpretations, and the distinction is context  
2 dependent. Therefore, it is possible to limit the argument to a certain set of  
3 inferences such as evaluation of task accomplishment, or generalization to a  
4 specific universe of observation without the necessity of theory. For example,  
5 if the target is the piano performance, then scores can be interpreted as  
6 observable attributes without a need to generalize beyond test. Therefore, for  
7 observable attributes, interpretive argument involves only the inferences of  
8 evaluation, generalization, and extrapolation. However, for theoretical  
9 constructs, one more inference is needed to explain scores in terms of a  
10 construct and to interpret them as indicators of specific abilities (Kane, et al.,  
11 1999).

### 13 **CURRENT TRENDS AND FUTURE DIRECTIONS**

15 Validity is currently defined in Standards parallel to Messick's construct  
16 validity model as "the degree to which evidence and theory support the  
17 interpretations of test scores entailed by proposed uses of tests" (AERA, APA,  
18 & NCME, 1999, p.9). However, Borsboom et al., (2004) claim that despite all  
19 the evolutions in the concept of validity over the years, when asked, most  
20 researchers in the field of psychology still define validity as "whether a test  
21 measures what one intends to measure" probably due to the failure of the  
22 construct validity model in providing a clear and workable conceptual scheme  
23 for practitioners.

24 In language performance assessment however, Kane's interpretative  
25 argument seems to be acknowledged to offer a feasible plan for validation,  
26 and has already been put into practice by Chapelle et al., (2004, in McNamara  
27 & Roever, 2006) in validating TOEFL. Recently elaborating on Kane's  
28 model, Bachman (2005) has developed the "assessment use argument" as a  
29 conceptual framework based on Toulmin's structure of reasoning involving  
30 claims, warrants, evidence, and rebuttals to achieve validation. Bachman's  
31 argument consists of two parts: 1) an assessment utilization argument that  
32 links assessment performance to a decision; and 2) an assessment validity  
33 argument that links the assessment performance to an interpretation.  
34 According to this model, since the aim is to justify a specific assessment, a  
35 "local theory" is sufficient to make claims about the decisions and  
36 interpretations based on the assessment, and to determine the types of  
37 evidence that needs to be collected to support these claims.

38 Although Bachman's model seems to be comprehensive and practical at  
39 first glance, given the context-dependent intricate interactions inherent in L2  
40 construct during a performance, it may still be a problem to require a theory  
41 base for all assessment types. For example, the social interactive view states  
42 that constructs in performances are co-constructed through social interactions,

## DEVELOPMENTS IN VALIDITY RESEARCH IN SECOND LANGUAGE PERFORMANCE TESTING

1 socially and culturally embedded, and context- dependent; therefore, ability,  
2 ability in language user, and context are inseparable making it impossible to  
3 measure the underlying abilities (Chalhoub-Deville, 2003).

4 Therefore, the distinction made by Kane between theoretical constructs  
5 and observable traits without requiring a theory all the time is a sensible  
6 approach. As Chalhoub-Deville & Deville (2005) suggest, according to the  
7 purpose and context, it is possible to seek different validity arguments and  
8 prioritize evidence for a particular use or decision-making. If we are  
9 interested just in the performance/task fulfillment, replicability and  
10 generalizability would not be the issue; however, if we are interested in  
11 performance with relation to the construct definition/language ability,  
12 generalizability would be necessary because the consistency or variability of  
13 performances contributes to the score meaning.

14 In conclusion, linking the performances to a theory to be able to interpret  
15 the results in terms of abilities and accordingly to be able to generalize is the  
16 most problematic area in validation. Therefore, further research is needed  
17 primarily to determine constructs which are very complex and elusive in  
18 performance assessment. Chalhoub-Deville (2003) suggests that it might be  
19 possible to determine any stable constructs that are accessed in similar ways  
20 across contexts by analyzing tasks and interacting factors in performance  
21 assessments in different contexts especially through ethno-methodological  
22 research so that the association networks used in varied situations to transfer  
23 knowledge and skill can be understood, and generalizability across contexts  
24 can be achieved. In addition, as social consequences of tests are not  
25 adequately integrated in validation models, development of a social theory  
26 regarding the social and political context in which assessment takes place  
27 should also be considered to understand the potential sources of unfairness  
28 and the meaning of test use in context (McNamara, 2006).

**BIBLIOGRAPHY**

- 1  
2  
3 American Psychological Association (APA) (1985). *Standards for*  
4 *educational and psychological testing*. Washington, DC.  
5 AERA, APA, & NCME (1999). *Standards for educational and*  
6 *psychological testing*. Washington, D.C.  
7 Bachman, L.F. (1990). *Fundamental considerations in language testing*.  
8 Oxford: Oxford University Press.  
9 Bachman, L. F. (2002, a) Alternative interpretations of alternative  
10 assessments: Some validity issues in educational performance assessments.  
11 *Educational Measurement: Issues and Practice*, 2(3), 5–18.  
12 Bachman, L. F. (2002b). Some reflections on task-based language  
13 performance assessment. *Language Testing*, 19, 453-476.  
14 Bachman, L. F. (2005). Building and supporting a case for test use.  
15 *Language Assessment Quarterly*, 2, 1-34.  
16 Borsboom, D., Mellenbergh, G. J. & Van Heerden, J. (2004). The concept  
17 of validity. *Psychological Review*, 111 (4), 1061-1071.  
18 Chalhoub-Deville, M. (2003). Second language interaction: current  
19 perspectives and future trends. *Language Testing*, 20 (4), 369-383.  
20 Chalhoub-Deville, M. & Deville, C. (2005). A look back at and forward to  
21 what language testers measure. In Hinkel, E. (ed.). *Handbook of research in*  
22 *second language teaching and learning* (pp. 815-831).Mahwah, NJ: Lawrence  
23 Erlbaum.  
24 Chapelle, C. A. (1999). Validity in language assessment. *Annual Review*  
25 *of Applied Linguistics*, 19, 254-272.  
26 Cronbach, L. J. & Meehl, P. E. (1955). Construct validity in psychological  
27 tests. *Psychological Bulletin*, 52, 281-302.  
28 Cronbach, L.J. (1988). Five perspectives on validation argument. In H.  
29 Wainer and H. Braun (eds.) *Test validity* (pp. 3-17). Hillsdale, NJ: L.Erlbaum.  
30 Kane, M., Crooks, T., Cohen, A. (1999). Validating measures of  
31 performance. *Educational Measurement: Issues and Practice*, 18 (2), 5-17  
32 Kane, M. (2001). Current concerns in validity theory. *Journal of*  
33 *Educational Measurement*, 38 (4), 319-342  
34 Lynch, B. K. (2001). Rethinking assessment from a critical perspective.  
35 *Language Testing*, 18 (4), 351-372  
36 McNamara, T. (1996). *Measuring second language performance*. London:  
37 Longman  
38 McNamara, T. (1997). “Interaction” in second language performance  
39 assessment: Whose performance? *Applied Linguistics*, 18 (4), 446-466.  
40 McNamara, T. & Roever, C. (2006). *Language Testing: Social dimension*.  
41 Oxford: Blackwell.

DEVELOPMENTS IN VALIDITY RESEARCH IN SECOND LANGUAGE  
PERFORMANCE TESTING

1       McNamara, T. (2007). Language assessment in foreign language  
2 education: the struggle over constructs. *The Modern Language Journal*, 91  
3 (2), 280-282.

4       Messick, S. (1989). Validity. In R. L. Linn (Ed.), *Educational*  
5 *measurement* (3<sup>rd</sup> ed., pp. 13-103). New York: Macmillan.

6       Messick, S. (1994). The interplay of evidence and consequences in the  
7 validity of performance assessment. *Educational Researcher*, 23: 2, 13-23

8       Messick, S. (1995). Validity of psychological assessment: Validation of  
9 inferences from persons' responses and performances as scientific inquiry into  
10 score meaning. *American Psychologist*, 50, 741-749.

11       Moss, P. A. (1994). Can there be validity without reliability? *Educational*  
12 *Researcher*, 23, 5-12.

13       Shepard, L. A. (1993). Evaluating test validity. *Review of Research in*  
14 *Education*, 19, 405-450.

15       Shohamy, E. (2001). *The power of tests: A critical perspective on the uses*  
16 *of language tests*. London: Longman

17