# RATING SCALES FOR SPEAKING TESTS: IDENTIFYING SOME CHALLENGES

**Ravshanova Ruzigul**

## Annotation

This paper aims to survey some of the problems associated with rating scales used in speaking tests. It will explore various types of rating scales, including analytic, holistic, and dynamic scales, highlighting their strengths and limitations. Furthermore, the paper will delve into the challenges of subjectivity, cultural bias, limited descriptors, and the overemphasis on language form, which can impact the reliability of speaking assessments. By identifying these issues, the paper seeks to contribute to the ongoing discourse on best practices in language testing and offer recommendations for enhancing rating scales.

**Key words:** rating scale, speaking proficiency, standardization, feedback, dynamic scales, descriptor, rater training, calibration.

### *The Role of Rating Scales*

Rating scales are systematic tools that facilitate the assessment of speaking performance by providing clear criteria against which candidates are evaluated. They offer several advantages in the context of speaking tests:

Standardization: Rating scales promote uniformity in assessment, ensuring that all candidates are judged by the same standards (Green, 2009). This standardization is essential for maintaining the reliability and validity of assessments, as it reduces the potential for rater bias and subjective interpretations.

Guidance for Raters: Scales provide explicit descriptions of performance levels, guiding raters in their evaluations. Clear criteria help raters make informed decisions about candidates' speaking abilities, leading to more objective scoring (Fulcher, 2003).

Feedback Mechanism: Effective rating scales not only assess performance but also offer valuable feedback for learners. By highlighting specific strengths and

weaknesses, scales enable candidates to understand their performance better and identify areas for improvement (Bygate, 2018).


*Challenges of Rating Scales*

Despite their widespread use and acknowledged benefits, rating scales for speaking tests are not without significant challenges. These issues can compromise the reliability, validity, and overall effectiveness of assessments. This section examines four major challenges associated with rating scales: subjectivity and bias, cultural differences, limited descriptors, and overemphasis on language form.

1. Subjectivity and Bias

One of the most prominent challenges in using rating scales is the inherent subjectivity associated with rater evaluations. Raters bring their interpretations, experiences, and biases into the assessment process, which can lead to inconsistencies in scoring (Bygate, 2018). For instance, two raters may interpret the same performance differently based on their individual standards and expectations. This subjectivity undermines the reliability of the scores, as candidates might receive vastly different ratings for similar performances.

Research has shown that raters' backgrounds, including their training, teaching experiences, and cultural contexts, can influence their scoring decisions (Green, 2009). To address this issue, it is essential to implement rater training programs that focus on calibration and consistency. By providing raters with clear guidelines and examples of performances at different levels, training can help minimize subjectivity and create a more standardized assessment environment.

2. Cultural Differences

Cultural differences pose another significant challenge to the use of rating scales in speaking assessments. Language use is heavily influenced by cultural norms, and raters from different cultural backgrounds may prioritize different aspects of communication (Green, 2009). For example, in some cultures, directness and

assertiveness in communication are valued, while in others, indirectness and politeness may be preferred. These cultural nuances can lead to disparities in how performances are rated.

The challenge lies in ensuring that rating scales are culturally inclusive and sensitive to diverse communicative norms. If a scale is biased toward a particular cultural perspective, it may unfairly disadvantage candidates from other backgrounds (Bygate, 2018). To mitigate this issue, it is crucial to involve diverse stakeholders in the development of rating scales, ensuring that multiple perspectives are considered. Additionally, raters should receive training on cultural competence to enhance their understanding of different communication styles and norms.

3. Limited Descriptors

Another challenge associated with rating scales is the potential for limited or vague descriptors within the scales themselves. When descriptors lack clarity or specificity, raters may struggle to accurately assess speaking performance (Fulcher, 2003). For example, a scale that states a candidate has "good fluency" without defining what constitutes "good" can lead to varied interpretations. Raters may apply their standards, resulting in inconsistent evaluations.

To enhance the effectiveness of rating scales, it is essential to develop detailed descriptors that clearly outline performance expectations for each level. Providing examples of specific language behaviors and communication strategies can help raters better understand what constitutes proficiency at different levels (Skehan, 1998). Comprehensive descriptors not only improve scoring reliability but also offer valuable feedback for candidates, enabling them to identify areas for improvement.

4. Overemphasis on Language Form

Many rating scales tend to place a disproportionate emphasis on linguistic accuracy, which can overshadow other critical aspects of speaking proficiency, such as interactional competence and communicative effectiveness (Skehan, 1998). This overemphasis may lead raters to prioritize grammatical correctness over the ability

to engage effectively in conversation, which is a key component of successful communication in real-world contexts.

For instance, a candidate who makes grammatical errors but demonstrates strong interactive skills may receive a lower score than warranted, while another candidate who is grammatically accurate but lacks fluency and engagement may receive a higher score (Fulcher, 2003). To address this challenge, rating scales should strive for a more balanced assessment that considers multiple dimensions of speaking performance, including fluency, interaction, and appropriateness of language use in context.

### *Recommendations for Improving Rating Scales*

Improving the effectiveness of rating scales for speaking assessments is essential for ensuring fair and accurate evaluations of language proficiency. Given the challenges discussed earlier, several strategies can be implemented to enhance the reliability, validity, and overall effectiveness of these assessment tools. This section outlines key recommendations that can contribute to the improvement of rating scales in speaking tests.

1. Develop Detailed and Clear Descriptors

One of the primary issues with many rating scales is the use of vague or limited descriptors that do not provide clear guidance for raters. To address this challenge, it is essential to develop detailed and specific descriptors that outline performance expectations at each scoring level. These descriptors should encompass various aspects of speaking proficiency, including fluency, accuracy, pronunciation, vocabulary usage, and interactional skills.

For instance, instead of a generic descriptor like "good fluency," a more detailed description might state, "speaks smoothly with minimal hesitation and can maintain a conversation without frequent pauses." Such specificity allows raters to understand better what constitutes proficiency at different levels and improves consistency in scoring (Fulcher, 2003). Moreover, examples of specific language

behaviors that illustrate each descriptor can further enhance clarity and provide raters with concrete benchmarks for assessment.

2. Implement Rater Training and Calibration

Rater training is a critical component of ensuring consistency and reliability in speaking assessments. Educators and raters should participate in comprehensive training programs that focus on the effective use of rating scales. This training should include sessions on understanding the descriptors, recognizing common pitfalls in scoring, and developing strategies to minimize subjectivity and bias.

Calibration sessions, where raters evaluate the same performances and discuss their scoring decisions, can be particularly effective in aligning expectations and interpretations (Green, 2009). Such collaborative efforts help raters develop a shared understanding of the assessment criteria, thereby increasing scoring reliability.

Regular refresher training sessions can also be beneficial to maintain rater consistency over time.

3. Foster Cultural Competence Among Raters

Given the impact of cultural differences on speaking assessments, it is crucial to foster cultural competence among raters. Raters should be trained to recognize and appreciate diverse communicative norms and styles, enabling them to evaluate performances fairly and inclusively. This training can include workshops on cultural awareness, discussions about the influence of cultural factors on communication, and strategies for avoiding cultural bias in assessments.

Additionally, involving diverse stakeholders in the development of rating scales can help ensure that multiple cultural perspectives are represented. By integrating a broader range of communicative practices and expectations into the rating scales, assessments can become more equitable and reflective of varied language use in different contexts (Bygate, 2018).

4. Balance Focus on Language Form and Communication

To address the overemphasis on language form in many rating scales, it is essential to develop scales that balance the evaluation of linguistic accuracy with other critical aspects of speaking proficiency, such as fluency and interactional competence. Rating scales should recognize that effective communication goes beyond grammatical correctness and encompasses the ability to engage in meaningful interactions.

For example, scales could incorporate criteria that assess candidates' ability to initiate and maintain conversations, respond appropriately to questions, and demonstrate active listening skills. By broadening the scope of evaluation, educators can create a more holistic assessment of speaking proficiency that reflects the complexities of real-world communication (Skehan, 1998).

**References:**

1. Bygate, M. (2018). Rating scales for speaking tests: A survey of some problems. Cambridge University Press.

2. Canale, M., & Swain, M. (1980). Theoretical bases of communicative approaches to second language teaching and testing. Applied Linguistics, 1(1), 1-47.

3. Fulcher, J. (2003). Testing second language speaking. Pearson Education.

4. Green, A. (2009). Language testing in practice. Oxford University Press.

5. Littlewood, W. (2004). The task-based approach: Some questions and suggestions. ELT Journal, 58(4), 319-325.

6. Skehan, P. (1998). A cognitive approach to language learning. Oxford University Press.