# A Driven Listening–Speaking Assessment Scale: Examining the Three-Dimensional Linkage of Usability-Acceptability-Practical Instructional Value!

## Li Gao & Ahmad Hishamuddin[*]
Faculty of Human Development, Sultan Idris Education University
Corresponding Author Email: hishamuddin.a@fpm.upsi.edu.my

**Abstract**
This study addresses limitations of the CED 4/6 listening–speaking scales, including insufficient subskill coverage, crude scoring items, and missing quantitative anchors and scenario exemplars. Based on psychometric reliability and validity theory, we designed a high-precision assessment scale that integrates usability → acceptability → practical instructional value. The scale decomposes speaking ability into eight subdimensions: pronunciation clarity, prosodic rhythm, lexical diversity, syntactic complexity, emotional expression, cross-cultural pragmatics, technical terminology mastery, and interaction strategies. Each subdimension includes a five-level quantitative rubric with positive and negative exemplars and scenario demonstration videos. This design reinforces the operability of scoring anchors. The system incorporates a dynamic proficiency-tiering mechanism. It matches instructional plans and tailored practice resources based on real-time scores. This creates a closed-loop assessment–feedback–intervention model. We assessed scale reliability and validity using Delphi consistency tests, KMO and Bartlett's tests, exploratory factor analysis, and Cronbach's α. In the empirical phase, we conducted Pearson correlation, paired-sample difference tests, and hierarchical regression path analysis. These analyses confirmed that interface usability and exemplar clarity have significant positive effects on teacher trust and continued use. They also demonstrated that acceptability mediates the effect of usability on practical instructional value. Results indicate that the new scale significantly enhances inter-rater consistency and teacher adoption. It supports personalized instruction in art education across regions. These findings offer a replicable methodological paradigm for the digital and intelligent transformation of educational assessment.
**Keywords:** Usability, Acceptability, Practical Instructional Value, Psychometrics, Listening–Speaking Assessment Scale

**Introduction**
Current CED 4/6 scales exhibit clear deficiencies in dimension coverage, item granularity, reliability, and regional applicability. For example, 42.3% of universities assess only fluency.

39.5% of scales focus solely on accuracy. Only 16.8% of scales address key subskills such as pronunciation clarity, prosodic rhythm, and pragmatic strategies. Item wording remains vague. The scales lack unified scoring anchors and situational examples. This inconsistency limits inter-rater reliability. The design does not capture lexical diversity, syntactic complexity, or cross-cultural pragmatic differences. Consequently, it cannot provide precise guidance for instructional improvement. A survey of Chengdu art students found that 65% of teachers feel the CED 4/6 scales inadequately support art students' linguistic characteristics. Moreover, the scales do not adapt to student characteristics or artistic expression habits. Their regional applicability remains low. Therefore, a novel speaking assessment scale must be developed. This scale should adopt a framework of usability→acceptability→practical instructional value. It must cover phoneme clarity, prosodic rhythm, specialized terminology usage, and interactive response strategies. It should include both positive and negative scoring examples and clear operational guidelines. The scale should allow flexible adjustments based on weighted student ability dimensions. Such a scale would help teachers objectively and accurately assess students' true speaking abilities. It would align assessment results with instructional practice.

Traditional listening–speaking scales rely on vague labels such as "natural expression," "smooth communication," and "accurate pronunciation." They omit quantitative benchmarks or scenario-based examples for critical subskills like lexical diversity, emotional expression, and discourse coherence. As a result, teachers cannot objectively differentiate performance across dimensions or deliver targeted feedback (Dvorski et al., 2022). Moreover, current scales do not include positive and negative exemplars or clear scoring anchors. This omission undermines scoring consistency and reproducibility (De, 2023). The development process also overlooks variations in student proficiency and disciplinary background. In particular, it fails to accommodate art students' unique phonetic traits and expression habits. Consequently, these scales do not meet art students' individual needs nor scale effectively to other regions (Sabrila & Apoko, 2022). To address these gaps, this study proposes a listening–speaking assessment scale with multiple innovations in both content and structure. First, core subskills are broken down into actionable tiers. Positive and negative examples accompany each tier to ensure clear scoring guidelines (Joseph, 2024). Second, the scale supports diverse assessment formats. It includes illustrated user guides and demonstration videos. Usability testing indicates the system is intuitive and lowers the learning curve substantially (Olagundoye et al., 2024). Third, the scale integrates Chengdu art students' prosodic patterns and proficiency profiles into parameterized weighting. Users can adjust settings based on student backgrounds. This feature enables seamless application across regions. Finally, to support students with weaker speaking skills, the system automatically suggests tailored practice resources based on assessment outcomes. It then generates personalized lesson plans targeting identified weaknesses. This closed-loop mechanism of assessment, feedback, and intervention promotes continuous speaking improvement.

Grounded in educational measurement reliability and validity research, this study integrates scale usability evaluation with psychometric rigor. It empirically validates a three-dimensional linkage model: usability → acceptability → practical instructional value. This approach transcends the traditional focus on scale internal consistency and construct validity. It establishes a framework where interface experience, teacher adoption behavior, and instructional outcomes reinforce one another (Vlachogianni & Tselios, 2022). In practice,

we developed a parameterized, multi-tier listening–speaking scale. It yields extensive empirical data via subskill tiering, positive and negative exemplars, and scenario anchors. The scale supports universities and education authorities in formulating adaptive assessment standards. It also enables flexible use across diverse disciplines and regions. Additionally, it assists teachers with tiered instruction and personalized tutoring (Aryadoust, 2023). Moreover, the scale's parameterized weights adjust flexibly based on student proficiency profiles. This feature enables seamless cross-regional adaptation. The research introduces innovations in multi-level quantitative scoring, integration of scenario exemplars with parameterized weights, and a closed-loop assessment–feedback–intervention mechanism. These advances address the lack of scales for evaluating cross-cultural pragmatics and interaction strategies (Davis et al., 2024). This study offers a systematic methodology and a replicable best-practice model for designing and evaluating listening–speaking assessment scales. By markedly improving scale usability and acceptability, it achieves deep integration with practical instructional value. These contributions lay a robust policy and practice foundation for the digitalization, intelligence, and personalization of educational assessment.

## Theoretical Background
### *High-Usability Scale Design and Teaching Value*

This study presents a listening–speaking assessment scale grounded in psychometric principles and scale usability evaluation. It balances rich content with ease of use (Scanferla et al., 2023). The scale specifies eight subskills: pronunciation clarity, prosodic rhythm, vocabulary usage, syntactic complexity, emotional expression, cross-cultural pragmatics, technical terminology mastery, and interactive response strategies. It establishes multi-tier quantitative scoring. Each tier includes positive and negative exemplars to help teachers decompose oral performance into discrete evaluation units (Mardon et al., 2025). The interface employs a modular layout with intelligent prompts and guided key steps. This reduces the learning curve and enhances evaluation consistency and efficiency. Teachers report that the scale's guidelines and exemplars are clear and easy to understand (Ofosu-Ampong, 2024). Grounded in educational psychology and implementation theory, the scale continuously incorporates teacher feedback on rubric clarity and exemplar relevance, leading to iterative refinements that significantly enhance teacher trust and willingness to use the tool (Rebollo & De, 2024). Furthermore, for students with lower English proficiency, assessment outcomes inform the design of targeted, personalized exercises or instructional plans, enabling teachers to tailor practice materials and teaching strategies to address each learner's specific needs. It then generates personalized tutoring plans. This closed-loop assessment–feedback–intervention model supports tiered instruction and individualized coaching. It effectively enhances student speaking skills and teaching quality.

### *Personalized and Proficiency-Tiered Assessment Matching*

This study decomposes the listening–speaking scale into eight subdimensions: pronunciation clarity, prosodic rhythm, lexical diversity, syntactic complexity, emotional expression, cross-cultural pragmatics, technical terminology usage, and interaction strategies. For each subdimension, it establishes a five-level scoring rubric with scenario-based exemplars. This design ensures thorough content coverage and user-friendly operation (Vlachogianni & Tselios, 2022). The system architecture incorporates a dynamic proficiency-tiering mechanism. It automatically maps subdimension scores to proficiency levels. It then selects appropriate test forms and practice materials (Toolaroud et al., 2023). Next, the

system generates targeted practice items and explanatory exemplars based on identified weaknesses. It also designs a sequenced learning path to form a personalized teaching plan. This enables precise instructional intervention (Moore et al., 2024). The interface employs a modular layout with guided workflows. This design markedly reduces the teacher learning curve and enhances assessment efficiency. Driven by highly acceptable feedback, teachers report greater trust in the tiered reports and a stronger intention to continue using the system. Ultimately, the tiered assessment outcomes directly inform group-based instruction and individualized tutoring. Students' performance improves significantly. This completes the closed loop of usability→acceptability→practical instructional value.

**Research Design and Methods**

*Research Design and Sampling*

We used a mixed paradigm centered on quantitative methods, supplemented by psychometric reliability and validity checks. We then built a three-stage path model: usability → acceptability → practical instructional value. In the preliminary phase, following Ma et al. (2023), we collected data on rubric clarity, exemplar interpretability, and inter-rater agreement. We conducted content coverage analyses and assessed structural and construct validity. We then drew on established educational measurement frameworks to define metrics for acceptability and practical instructional value (Musa et al., 2024). We used three-dimensional stratified random sampling based on faculty rank, teaching experience, and institution type. Each stratum included at least ten teachers. There were 18 strata, representing 180 potential participants. A priori power analysis with G*Power 3.1 ($f^2 = 0.15$, $\alpha = 0.05$, power = 0.80) indicated a minimum sample size of 77. After accounting for invalid cases and multiple-group comparisons, we aimed to secure at least 50 valid responses. Data were collected using the Tencent Questionnaire platform. We implemented logic branching, required responses, and an informed consent form. The survey link was shared in WeChat groups with multiple reminders. We then screened out responses with abnormal durations, uniform answers, or excessive missing data. We retained 50 high-quality responses for reliability, validity, and regression analyses.

*Development and Implementation of the Measurement Scale*

We developed and implemented the measurement scale in three stages. First, we conducted a systematic literature review and semi-structured interviews. We used a construct-validity and scale usability framework to identify potential dimensions of listening–speaking assessment administration and evaluation workflows (Murray-Smith et al., 2022). We then gathered core items—trust, adoption intent, report usability, and student gains—from 10 teachers and 3 instructional designers. Next, we conducted three Delphi rounds with five experts in assessment and applied linguistics. They provided anonymous ratings and feedback. After revisions, the third round achieved a Kendall's W > 0.75 (p < .01), confirming the item structure. We then organized the survey into three subscales—usability, acceptability, and teaching practical value. It comprised 30 five-point Likert items, including reversed items to control bias. Finally, we piloted the draft survey with 30 teachers. We ran exploratory factor analysis (principal components, eigenvalue > 1, Varimax rotation), computed Cronbach's α, and checked item-total correlations (r > .30). We removed items with low loadings or correlations. We also refined and merged others. The final scale had a clear structure and strong internal consistency. It laid the groundwork for large-scale reliability, validity, and regression analyses.

*Data Analysis Strategy*

We first computed Cronbach's α for the full scale and each subscale. We examined item-deleted α changes and removed low-contributing items to ensure internal consistency. Next, we assessed sampling adequacy with the Kaiser–Meyer–Olkin measure (KMO > .60) and sphericity with Bartlett's test (p < .001). Only then did we proceed with factor analysis. After reliability and validity testing, we examined how teacher background influenced dimension scores. We first ran Pearson correlations among usability, acceptability, and instructional practical value. Next, we used paired-sample t-tests to compare adjacent dimension means. If normality or homogeneity assumptions failed, we applied nonparametric tests or Welch correction to ensure robustness. Finally, we conducted hierarchical regression to test direct effects along the usability → acceptability → practical instructional value path. We entered usability in the first step. We then added acceptability and covariates. We reported coefficients and p-values for each predictor. We also checked multicollinearity by computing variance inflation factors (VIF < 5). If needed, we combined variables or applied principal component analysis. These steps ensured model validity and robustness.

**Empirical Analysis and Results**

*Reliability and Validity Assessment*

We first assessed internal consistency using Cronbach's α. Table 1 shows the full scale (20 items) α = .819 (standardized = .821), exceeding the .80 threshold. Usability (6 items) registered α = .683, and practical value (8 items) α = .687—both within acceptable range. Acceptability (6 items) had α = .610, which is marginally low. This suggests reviewing item–total correlations and revising those items.

Table 1
*Reliability Statistics*

| Scale or Subscale | Cronbach's α | Standardized Items Cronbach's α If Items Are Standardized | Number of Items |
|---|---|---|---|
| Overall | .819 | .821 | 20 |
| Usability | .683 | .689 | 6 |
| Acceptability | .610 | .592 | 6 |
| Practical Value | .687 | .698 | 8 |

We then assessed data suitability for factor analysis. We used the Kaiser–Meyer–Olkin measure and Bartlett's test of sphericity. Table 2 reports KMO and Bartlett statistics for each subscale. The usability subscale had KMO = .598 and $\chi^2(15) = 34.05$, p = .003. The acceptability subscale had KMO = .622 and $\chi^2(15) = 36.19$, p = .002. The practical value subscale had KMO = .597 and $\chi^2(28) = 77.14$, p < .001. Although KMO values fell just below the .70 benchmark, significant Bartlett tests confirmed sufficient inter-item correlation for factor extraction.

Table 2
*KMO Measure and Bartlett's Test of Sphericity*

| Usability | | | Acceptability | | | Practical Value | | |
|---|---|---|---|---|---|---|---|---|
| KMO Measure of Sampling Adequacy | | .598 | KMO Measure of Sampling Adequacy | | .622 | KMO Measure of Sampling Adequacy | | .597 |
| Bartlett's Test of Sphericity | Approx. Chi-Square | 34.052 | Bartlett's Test of Sphericity | Approx. Chi-Square | 36.189 | Bartlett's Test of Sphericity | Approx. Chi-Square | 77.136 |
| | Df | 15 | | Df | 15 | | Df | 28 |
| | Sig. | .003 | | Sig. | .002 | | Sig. | .000 |

Based on these findings, we will proceed to exploratory and confirmatory factor analyses to evaluate construct validity. High α coefficients and significant Bartlett results provide robust evidence of scale reliability and structural validity. They also lay a firm foundation for subsequent mean difference tests and regression analyses.

*Dimension-Level Score Difference Tests*

We first ran Pearson correlations among usability, acceptability, and instructional practical value (Table 3). The correlations were as follows: usability and acceptability, $r = .289$, $p < .05$; usability and instructional practical value, $r = .365$, $p < .01$; acceptability and instructional practical value, $r = .460$, $p < .01$. These positive correlations suggest that interface usability, system trust, and perceived instructional support reinforce each other. Together, they underpin scale adoption and value recognition.

Table 3
*Correlation Coefficients*

| | | Usability | Acceptability | Practical Value |
|---|---|---|---|---|
| Usability | Correlation Coefficient | 1.000 | | |
| Acceptability | Correlation Coefficient | .289[*] | 1.000 | |
| Practical Value | Correlation Coefficient | .365[**] | .460[**] | 1.000 |

* Correlation is significant at the 0.05 level (two-tailed).
** Correlation is significant at the 0.01 level (two-tailed).

Next, we compared mean scores for each dimension within the same teacher sample (n = 30) (Table 4). The difference between usability and acceptability was not significant ($\Delta M = -0.067$, SD = 2.651, $t = -0.138$, $p = .891$). However, usability and instructional practical value differed significantly ($\Delta M = -8.633$, SD = 2.632, $t = -17.963$, $p < .001$, 95% CI [−9.616, −7.650]). Similarly, acceptability and instructional practical value differed significantly ($\Delta M = -8.567$, SD = 2.402, $t = -19.532$, $p < .001$, 95% CI [−9.464, −7.670]). These results indicate that teachers rated instructional support significantly higher than interface usability and overall acceptability.

Table 4
*Paired Samples Test*

| | Mean Difference | | | | | | | |
| | Mean | Std. Deviation | Std. Error Mean | 95% Confidence Interval of the Difference | | t | Df | Sig. (2-tailed) |
| | | | | Lower | Upper | | | |
|---|---|---|---|---|---|---|---|---|
| Usability - Acceptability | -.067 | 2.651 | .484 | -1.057 | .923 | -.138 | 29 | .891 |
| Usability - Practical Value | -8.633 | 2.632 | .481 | -9.616 | -7.650 | -17.963 | 29 | .000 |
| Acceptability - Practical Value | -8.567 | 2.402 | .439 | -9.464 | -7.670 | -19.532 | 29 | .000 |

In sum, correlations and paired-sample tests reveal that, despite intercorrelations, teachers perceived instructional practical value as significantly higher than usability and acceptability. These findings suggest that future listening–speaking assessment scale optimizations should prioritize instructional support and feedback. Simultaneously, interface usability and system trust should be enhanced. This approach will align functionality with user experience.

*Regression Analysis*

We modeled instructional practical value as the dependent variable. Usability and acceptability served as predictors in a multiple regression. Table 5 shows that the model outperforms the null model, $F(2, 27) = 10.145$, $p = .001$. The $R^2$ of .429 indicates that predictors explain 42.9% of variance.

Table 5
*ANOVA[a]*

| Model | | Sum of Squares | Df | Mean Square | F | Sig. |
|---|---|---|---|---|---|---|
| 1 | Regression | 99.028 | 2 | 49.514 | 10.145 | .001[b] |
| | Residual | 131.772 | 27 | 4.880 | | |
| | Total | 230.800 | 29 | | | |

a. Dependent Variable:Practical Value
b. Predictors: (Constant),Acceptability, Usability

As shown in Table 6, usability predicts instructional practical value (B = 0.407, β = .334, p = .039). Acceptability also predicts value (B = 0.572, β = .462, p = .006). Both effects are significant, with acceptability slightly stronger.

Table 6
*Regression Coefficientsa[a]*

| Model | | Unstandardized Coefficients B | Std. Error | Standardized Coefficients Beta | t | Sig. | 95% Confidence Interval for B Lower | Upper | Correlations Zero-order | Partial | Part | Collinearity Statistics Tolerance | VIF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Constant | 9.146 | 5.828 | | 1.569 | .128 | -2.811 | 21.103 | | | | | |
| | Usability | .407 | .188 | .334 | 2.166 | .039 | .022 | .793 | .489 | .385 | .315 | .888 | 1.126 |
| | Acceptability | .572 | .191 | .462 | 2.997 | .006 | .180 | .963 | .574 | .500 | .436 | .888 | 1.126 |

a. Dependent Variable:Practical Value

We checked multicollinearity by computing variance inflation factors and tolerances. Table 7 reports VIF = 1.126 and tolerance = .888 for both predictors. These values meet the VIF<5 and tolerance>0.2 standards. Although the maximum condition index reached 30.615, variance proportions were not focused on a single predictor. Thus, collinearity risk remains acceptable.

Table 7
*Collinearity Diagnostics[a]*

| Model | Dimension | Eigenvalue | Condition Index | Variance Proportions Constant | Usability | Acceptability |
|---|---|---|---|---|---|---|
| 1 | 1 | 2.992 | 1.000 | .00 | .00 | .00 |
| | 2 | .005 | 25.042 | .00 | .71 | .62 |
| | 3 | .003 | 30.615 | 1.00 | .29 | .38 |

a. Dependent Variable:Practical Value

In sum, regression results confirm that usability and acceptability independently predict instructional practical value. The model fits well, and collinearity is acceptable.

**Discussion**

This study used reliability and validity checks, score-difference tests, and regression analysis to validate the linkage among usability, acceptability, and practical instructional value. It also highlighted the central role of teacher experience. First, Cronbach's α and KMO/Bartlett tests confirmed good internal consistency and suitable factor structure. However, the acceptability subscale's reliability was slightly low. This suggests refining items on trust and continued use. Next, correlations were significantly positive across the three dimensions. Paired-sample t-tests showed that instructional practical value scores were significantly higher than usability and acceptability. This indicates that instructional support is key to system acceptance. Then, regression results showed that acceptability predicted instructional value more strongly than usability. VIF values were all below 5, and no serious collinearity was found. These findings support the usability → acceptability → practical instructional value pathway.

Theoretically, this study integrates scale usability theory with educational measurement validity models. This enriches research paradigms in technology acceptance and assessment validity. Practically, this work offers two recommendations for administrators and developers.

First, strengthen instructional feedback and decision support. Second, refine interface interaction and trust mechanisms. This study has limitations. The sample was drawn from undergraduate institutions in a single city and relied on teacher self-reports. Future work should include universities across regions and levels. It should incorporate multimodal data, such as student outcomes and classroom observations. Overall, this study provides empirical guidance for optimizing and scaling listening–speaking assessment scales. It advances cross-disciplinary integration and innovation in educational technology and assessment.

**Conclusion**

After identifying gaps in the CED 4/6 scales—dimension coverage, item granularity, scoring anchors, and regional applicability—this study proposes a new listening–speaking assessment scale. It centers on a three-dimensional linkage: usability→acceptability→practical instructional value. The scale decomposes speaking ability into eight subdimensions. Each subdimension includes positive and negative exemplars and scenario demonstrations. This setup enables multi-tier quantitative scoring. The system architecture incorporates student proficiency tiering and a closed-loop feedback mechanism. Based on assessment results, it automatically recommends personalized practice materials and instructional plans. We ensured scale reliability and validity by conducting Delphi consensus tests, exploratory factor analysis, and internal consistency checks. We also employed a modular interface design. Usability was evaluated through task completion time, error rates, and user satisfaction. Next, using the Technology Acceptance Model and hierarchical regression path analysis, we confirmed that interface usability and exemplar clarity significantly influence teacher trust and intention to continue use. We also highlighted how teaching feedback features promote effective classroom application. By integrating usability evaluation with psychometric theory, this study expands paradigms in technology adoption and assessment validity. It offers a replicable methodology for designing listening–speaking scales across regions and disciplines.

Based on these findings, we recommend enhancing the feedback architecture by enriching exemplar libraries and embedding real-time guidance prompts, refining the acceptability dimension through additional user-centered reviews to strengthen its consistency, broadening the empirical base by involving educators from varied regions and incorporating direct student performance metrics, developing comprehensive training modules and interactive tutorials to ensure consistent practitioner implementation, iterating scenario exemplars and quantitative anchors using longitudinal classroom data and end-user feedback, and exploring integration with complementary multimodal assessment technologies—such as speech analytics and video analysis—to drive the next generation of personalized, intelligent educational assessment.

## References

Dvorski Lacković, I., Kurnoga, N., & Miloš Sprčić, D. (2022). Three-factor model of Enterprise Risk Management implementation: exploratory study of non-financial companies. Risk Management, 24(2), 101-122.

De Jong, N. H. (2023). Assessing second language speaking proficiency. Annual Review of Linguistics, 9(1), 541-560.

Sabrila, R. A. P., & Apoko, T. W. (2022). The Effectiveness of Podcast on Listening Skill for Vocational School Students. IDEAS: Journal on English Language Teaching and Learning, Linguistics and Literature, 10(2), 1177-1186.

Joseph, R. (2024). Construct validity and reliability of the theory evaluation scale: A factor analysis. Journal of Social Work Education, 60(3), 295-309.

Olagundoye, O., Gibson, W., & Wagg, A. (2024). A protocol for the co-creation and usability/acceptability testing of an evidence-based, patient-centred intervention for self-management of urinary incontinence in older men. Plos one, 19(8), e0306080.

Vlachogianni, P., & Tselios, N. (2022). Perceived usability evaluation of educational technology using the System Usability Scale (SUS): A systematic review. Journal of Research on Technology in Education, 54(3), 392-409.

Aryadoust, V. (2023). The vexing problem of validity and the future of second language assessment. Language Testing, 40(1), 8-14.

Davis, F. D., Granić, A., & Marangunić, N. (2024). The technology acceptance model: 30 years of TAM. Switzerland: Springer International Publishing AG.

Scanferla, W. H., Oliveira, C., Lousada, M. L., & Teixeira, L. C. (2023). The Usability and Acceptability of the mHealth "Health and Voice" for Promoting Teachers' Vocal Health. Journal of Voice.

Mardon, A. K., Wilson, D., Leake, H. B., Harvie, D., Andrade, A., Chalmers, K. J., & Moseley, G. L. (2025). The acceptability, feasibility, and usability of a virtual reality pain education and rehabilitation program for Veterans: a mixed-methods study. Frontiers in Pain Research, 6, 1535915.

Ofosu-Ampong, K. (2024). Beyond the hype: exploring faculty perceptions and acceptability of AI in teaching practices. Discover Education, 3(1), 38.

Balanyà Rebollo, J., & De Oliveira, J. M. (2024). Teachers' evaluation of the usability of a self-assessment tool for mobile learning integration in the classroom. Education Sciences, 14(1), 1.

Vlachogianni, P., & Tselios, N. (2022). Perceived usability evaluation of educational technology using the System Usability Scale (SUS): A systematic review. Journal of Research on Technology in Education, 54(3), 392-409.

Toolaroud, P. B., Nabovati, E., Mobayen, M., Akbari, H., Feizkhah, A., Farrahi, R., & Jeddi, F. R. (2023). Design and usability evaluation of a mobile-based-self-management application for caregivers of children with severe burns. International wound journal, 20(7), 2571-2581.

Moore, S., Costello, E., Nguyen, H. A., & Stamper, J. (2024, July). An automatic question usability evaluation toolkit. In International Conference on Artificial Intelligence in Education (pp. 31-46). Cham: Springer Nature Switzerland.

Ma, B., Liu, Q., Jiang, Z., Che, D., Qiu, K., & Shang, X. (2023). Energy-Efficient 3D Path Planning for Complex Field Scenes Using the Digital Model with Landcover and Terrain. Isprs International Journal of Geo-Information, 12(2), 82.

Musa, H. G., Fatmawati, I., Nuryakin, N., & Suyanto, M. (2024). Marketing research trends using technology acceptance model (TAM): A comprehensive review of researches (2002–2022). Cogent business & management, 11(1), 2329375.

Murray-Smith, R., Oulasvirta, A., Howes, A., Müller, J., Ikkala, A., Bachinski, M., & Klar, M. (2022). What simulation can do for HCI research. Interactions, 29(6), 48-53.