# APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING IN PHYSICS EDUCATION

**Sujiyani Kassiavera,
A. Suparmi,
C. Cari,
Sukarmin Sukarmin**

**Abstract.** *The challenge of accurately assessing critical thinking in physics education, particularly on topics like work and energy, remains a key issue for educators. The current study aims to address this challenge by exploring students' critical thinking abilities using two-tier test data analyzed through the Rasch model. Data were collected from students in Bengkulu Province, Sumatra, Indonesia, and analyzed to evaluate item fit, reliability, and students' abilities across various critical thinking dimensions. It was found that the two-tier instrument demonstrated high validity and reliability, with infit and outfit mean square values close to ideal, and strong separation reliability for both participants and items. Further analysis revealed significant variations in students' abilities in aspects of critical thinking, including interpretation, analysis, and self-regulation, indicating the need for more targeted pedagogical interventions. The study concludes that applying the Rasch model to analyze two-tier tests not only enhances understanding of students' critical thinking but also provides a novel approach to developing and implementing evaluation instruments in physics education. These findings contribute to the existing literature by deepening theoretical insights into critical thinking within physics education and offering practical guidance for educators aiming to improve curriculum design and teaching strategies.*
**Keywords:** *critical thinking, item analysis, physics education, rasch model, reliability test, two-tier test*

**Sujiyani Kassiavera, A. Suparmi,
C. Cari, Sukarmin Sukarmin**
*Sebelas Maret University, Indonesia*

## Introduction

Critical thinking skills are crucial in physics education because they help students understand and apply basic concepts in complex real situations. A deep understanding of concepts such as work and energy cannot be achieved only through memorization; students need the skills to analyze, evaluate, and synthesize relevant information (Paul & Elder, 2014; ten Dam & Volman, 2004). Therefore, when solving physics problems, students are expected to identify relevant information, understand the assumptions underlying an argument, and develop logical and innovative solutions. Ennis (2018) and Facione (2020) have described critical thinking as a cognitive skill that allows individuals to analyze, evaluate, and synthesize information objectively and rationally. Moreover, it is essential in physics, where many problems require not only critical thinking but also a deep understanding of the relationships between various physical variables (Dessie et al., 2024; Wan, 2023). For example, in solving physics problems, students need to remember formulas or facts, understand the underlying principles, and apply them in real-life scenarios (Bao & Koenig, 2019; Nyirahabimana et al., 2023).

The evaluation of critical thinking in physics education is an essential component responsible for enhancing the quality of learning (van Laar et al., 2020; Wang et al., 2023). Evaluating critical thinking allows educators to measure the extent to which students can critically integrate and apply their knowledge, rather than simply recalling information (Memduhoğlu & Keleş, 2016; Sasson et al., 2018). Proper evaluation of these skills enables students to identify their strengths and weaknesses in understanding physics concepts (Dwyer & Walsh, 2020; Rapti & Sapounidis, 2024). Effective and structured evaluation instruments are necessary to accurately measure these skills (Mafinejad et al., 2017; Memduhoğlu & Keleş, 2016). Evaluation instruments should be designed to test students' ability to analyze, evaluate, and solve complex physics problems with a critical approach (Peng, 2023; Wang & Zhang, 2024). A well-designed evaluation should assess not only the final result but also the thinking process students use, thus providing a more comprehensive overview of their critical thinking (García-Carmona, 2023).

In recent years, two-tier tests have become a widely used type of evaluation designed to measure students' critical thinking in more detail, particularly in science education, including physics (Irmak et al., 2023; Zakwandi et

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/

ISSN 2538–7138 /Online/

al., 2024). Two-tier tests consist of two interrelated tiers (Affandy et al., 2024; Cetin-Dindar & Geban, 2011; Potvin et al., 2015). The first tier typically includes multiple-choice questions that test students' factual or conceptual knowledge, while the second tier asks students to provide reasons or explanations for the answers they selected in the first tier. These tests not only evaluate whether answers are correct or incorrect but also explore students' deeper understanding and their ability to explain the thought process behind their responses (Irmak et al., 2023; Zakwandi et al., 2024).

The relevance of two-tier tests in measuring critical thinking lies in their ability to determine whether students deeply understand concepts or merely memorize information without engaging in critical thinking (Affandy et al., 2024; Potvin et al., 2015). Two-tier tests require students to provide reasons or explanations, thereby revealing whether they can analyze, evaluate, and apply concepts in various situations—core elements of critical thinking (Affandy et al., 2024; Putica, 2023). Previous research has highlighted the effectiveness of two-tier tests in science education (e.g.: (Cari et al., 2020; Kaltakci et al., 2016; Treagust, 1988)). Using two-tier tests in physics education offers educators a clearer understanding of how students think critically and grasp complex concepts (Kaltakci et al., 2016; Potvin et al., 2015).

The utilization of the Rasch model in analyzing two-tier test data is essential because it provides a robust approach to measuring and understanding students' critical thinking skills in greater depth (Cvenic et al., 2022), enables data analysis that considers the difficulty level of each item as well as the ability of the individual being tested (Wang & Ho, 2024), providing a more accurate assessment of whether students not only know the correct answer (Cascella et al., 2020), but also are able to provide appropriate reasoning, which is an indicator of critical thinking ability (Laliyo et al., 2022). Using the Rasch model to measure critical thinking helps ensure that the test consistently evaluates students' critical thinking ability across different ability levels, thereby confirming that the assessment genuinely reflects the intended skills (Kaur et al., 2024).

### Research Problem

The primary problem addressed in this research is the inability to accurately measure critical thinking in physics education, an issue that has not received sufficient attention in previous studies. The measurement of critical thinking is often limited to evaluations that assess only factual knowledge or simple problem-solving skills, without delving into the analytical and evaluative thinking processes that underlie the understanding of complex physics concepts. This is a significant problem because critical thinking is essential to understanding, applying, and evaluating physics concepts in diverse and authentic contexts. The importance of accurately measuring critical thinking in physics education lies in its ability to equip students with the thinking skills needed to face future intellectual and professional challenges. If critical thinking skills are effectively measured and developed, physics education can focus on developing higher-order thinking skills, ultimately increasing the overall quality of education.

The inability to properly measure critical thinking can lead to misunderstandings about the extent to which students genuinely understand the material and result in unfair and superficial assessments of student competencies. Most evaluation methods used today are insufficiently sensitive to detect differences in the level of critical thinking among students, leading to less accurate assessment results and failing to provide a comprehensive picture of students' critical thinking. This gap highlights the need to develop more appropriate and structured evaluation instruments, such as two-tier tests analyzed with the Rasch Model, to ensure that critical thinking is measured validly, reliably, and thoroughly in physics education.

### Research Focus

Measuring critical thinking in physics education is essential as it significantly impacts the enhancement of educational quality and the development of theory in physics education. The utilization of the Rasch Model in the current study provides a more appropriate approach to evaluating student ability, as the Rasch Model allows for a detailed analysis of test items, including the identification of items that do not function well or fail to effectively differentiate between different levels of student ability. The positive impact of this research on the practice of physics education can be observed in several aspects. First, with improved evaluation instruments, educators can design teaching strategies that are more effective and tailored to the needs of students, thereby enhancing the overall quality of learning. Second, the study's findings can influence educational policy by encouraging the adop-

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

tion of more accurate methods to measure critical thinking skills, which are among the key competencies to be developed in 21st-century education. The research focus not only contributes to enhancing the quality of physics education but also supports the development of students who are better prepared to face future academic and professional challenges.

*Research Aim and Research Questions*

This study aimed to explore students' critical thinking skills using data from a two-tier test that was analyzed through the Rasch model. The exploration conducted in the research was intended to find a more accurate and in-depth evaluation method to measure critical thinking skills in the context of physics education, with the goal of enhancing the quality of teaching and learning. The study had several specific goals, including (1) evaluating the item suitability and reliability of the two-tier test used to measure students' critical thinking skills, (2) analyzing how well the two-tier test identified variations in students' critical thinking skills, including aspects such as interpretation, analysis, and self-regulation, and (3) identifying opportunities for more effective pedagogical interventions based on the research findings. The objectives of this research included: (1) collecting data from students to test the validity and reliability of the two-tier instrument, (2) analyzing the data using the Rasch model to measure infit and outfit fit, as well as separation reliability for participants and items, and (3) evaluating differences in students' critical thinking skills across the various dimensions tested to provide recommendations for the development of more targeted and effective teaching strategies. The research questions addressed in this study were: (1) How reliable were the two-tier test instruments in measuring students' critical thinking skills on the concept of work and energy in physics education? (2) How did Rasch model analysis help evaluate item fit on two-tier test instruments? and (3) What variations in students' critical thinking skills were identified through the Rasch model analysis of critical thinking dimensions such as interpretation, analysis, and self-regulation?

## Research Methodology

*General Background*

The current study used the Rasch model to evaluate students' critical thinking skills in physics education, particularly on the topic of work and energy. The research employed a quantitative design with a survey approach, where data were collected through a two-tier test and analyzed using the Rasch model. The study was conducted over three months, from June to August 2023, in Bengkulu province, Sumatra, Indonesia. The sample selection was carried out to obtain a broad representation of the diverse academic characteristics in Bengkulu province, including gender, age, semester, and study program. The study is based on the theory of critical thinking outlined by Facione, (2020) and Paul and Elder, (2014), which emphasizes the importance of the ability to analyze, evaluate, and synthesize information objectively in physics education. The Rasch model was chosen as the main analysis tool because it provides a deep understanding of the validity and reliability of measurement instruments in an educational context. The items were scored using the Partial Credit Model (PCM) method with four scoring categories by Affandy et al., (2021) and Istiyono, (2016). Data analysis included reliability measurement, item difficulty, item fit, and point measure correlation to evaluate the consistency and contribution of each item to the measurement of critical thinking skills.

*Participants*

The study population comprised 342 undergraduate students from various university programs in Bengkulu Province, Sumatra, Indonesia. The sample was drawn using a purposive sampling method, where the students selected were those involved in study programs relevant to physics topics, particularly the concepts of work and energy. Furthermore, exclusion criteria were applied, i.e., students who only attended up to 80% of the overall lectures were excluded from the sample. According to the statistical analysis with a margin of error of 5% and a confidence level of 95%, the minimum sample size required was 183 undergraduate students from various study programs.

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

**Table 1**
*Demographics of Respondents*

| Characteristic | Sub-characteristic | *N* | Percentage |
|---|---|---|---|
| Gender | Male | 84 | 45.9 |
| | Female | 99 | 54.1 |
| | Total | 183 | 100 |
| Age | 18 Years | 30 | 16.4 |
| | 19 Years | 22 | 12 |
| | 20 Years | 65 | 35.5 |
| | 21 Years | 66 | 36.1 |
| Semester | 1st semester | 30 | 16.4 |
| | 3rd semester | 33 | 18 |
| | 5th semester | 55 | 30.1 |
| | 7th semester | 65 | 35.5 |
| Department | Physics Education | 41 | 22.4 |
| | Science Education | 40 | 21.9 |
| | Civil Engineering | 30 | 16.4 |
| | Electrical Engineering | 32 | 17.5 |
| | Mechanical Engineering | 40 | 21.9 |

The study sample totaled 183 students, with 84 male students (45.9%) and 99 female students (54.1%). The age of the research participants varied from 18 to 21 years old, with 30 students aged 18 years old (16.4%), 22 students aged 19 years old (12.0%), 65 students aged 20 years old (35.5%), and 66 students aged 21 years old (36.1%). The sample was also distributed across various semesters, with 30 students in semester 1 (16.4%), 33 students in semester 3 (18.0%), 55 students in semester 5 (30.1%), and 65 students in semester 7 (35.5%). By study program, the participants included 41 students from the Physics Education study program (22.4%), 40 from Science Education (21.9%), 30 from Civil Engineering (16.4%), 32 from Electrical Engineering (17.5%), and 40 from Mechanical Engineering (21.9%).

The demographic data, presented in Table 1, provides an overview of the diverse characteristics of the participants, including gender, age, semester, and study program. These indicators are shown to illustrate the sample's representation of the broader population and to highlight the sample's relevance to the research objectives, which focus on evaluating critical thinking skills in physics education. Additionally, the demographic information offers a foundation for analyzing how these factors may influence the research results and gives deeper insight into the educational context in Bengkulu province. Participating students are guaranteed confidentiality by not disclosing personal information when publishing research results. It was done to ensure anonymity and voluntariness in the research. Furthermore, all data collected was processed anonymously, and any participation in the study was voluntary. Before data collection, undergraduate students were given a clear explanation of the purpose of the research and their right to withdraw from the study at any time without consequence.

*Instruments*

The instrument used was a two-tier test consisting of 20 items designed to measure students' critical thinking in an introductory physics course on work and energy. According to Facione, (2020), each item in the instrument addressed specific aspects of critical thinking skills: Interpretation (4 items), Analysis (3 items), Evaluation (4 items), Explanation (3 items), Inference (3 items), and Self-Regulation (3 items). The items were scored using the 4-category Partial Credit Model (PCM) method, adapted from the research results of Istiyono et al., (2019); Lukman et al., (2021); Zakwandi et al., (2024). The PCM scoring was organized into four categories: Category 1 was assigned when a participant provided both an incorrect multiple-choice answer and an incorrect reason, indicating a lack

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/

ISSN 2538–7138 /Online/

of understanding of the concept. Category 2 was given when the participant chose the correct multiple-choice answer but gave an incorrect reason, suggesting partial knowledge or guessing but weak conceptual understanding. Category 3 was applied when participants provided incorrect multiple-choice answers but correct reasoning, showing better knowledge even though they could not apply it correctly to the options. Finally, Category 4 was given if participants answered both the multiple-choice and reasoning questions correctly, reflecting full understanding and strong critical thinking skills.

*Data Analysis*

The first step in data analysis was to measure the reliability of the test, aimed at assessing the instrument's internal consistency. Reliability was calculated using the Rasch reliability coefficient, which provided an indication of how well the items, as a whole, measured critical thinking skills. After testing reliability, the analysis proceeded with an assessment of item difficulty (item measure). The analysis determined how difficult each item was for participants and ensured that the items covered an appropriate range of difficulty levels to measure critical thinking skills across different abilities.

The next step was an item fit analysis, which involved testing the Outfit Mean Square (MNSQ) and Outfit Z-Standard (ZSTD). MNSQ is a statistic that shows how well the data for each item fit the expected Rasch model, with an ideal value around 1.00. ZSTD provided additional information by identifying whether deviations from the model were statistically significant. Extremely high or low ZSTD values suggested that an item might not be functioning well in the context of the instrument. Finally, point measure correlation (pt mean corr) values were analyzed to determine how much each item contributed to the overall measurement of critical thinking skills. A positive and significant correlation indicated that the item was consistent with the overall measurement objectives.

## Research Results

*Overall Calibration of The Two-Tier Test*

The primary results from the calibration of the two-tier test measuring critical thinking on physics topics, specifically work and energy, indicated that the test items had fair measurement quality (see Table 2). The statistics showed that the standard error for participants was .34 logits, with a standard deviation of .03, while for items, the standard error was .14 logits, with a standard deviation of .01. The average infit MNSQ value for participants was .92 ($SD = 0.15$), and for items, it was .94 ($SD = 0.17$), suggesting that most items had a good level of fit with the Rasch model. The mean outfit MNSQ values also showed favorable results, with .97 ($SD = 0.21$) for participants and .99 ($SD = 0.19$) for items, indicating no significant deviations from the model.

The separation reliability for participants was .84, and for items, it was .94, demonstrating that the test had a reasonably strong ability to distinguish between participants' abilities and identify items of varying difficulty. The chi-square results were significant ($p < .05$) for both participants and items, confirming that the model could explain the variability in the data. Based on the overall calibration results, the test items were effective in measuring critical thinking on the topic of work and energy in physics, with a high degree of consistency and accuracy.

**Table 2**
*Summary Statistics from The Politomus Rasch Model*

| Statistic | Student ($N = 183$) | Item ($N = 20$) |
|---|---|---|
| Standard error (logit) | | |
| $M$ | 0.34 | 0.14 |
| $SD$ | 0.03 | 0.01 |
| Infit MNSQ | | |
| $M$ | 0.92 | 0.94 |
| $SD$ | 0.15 | 0.17 |
| Outfit MNSQ | | |
| $M$ | 0.97 | 0.99 |
| $SD$ | 0.21 | 0.19 |

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/

ISSN 2538–7138 /Online/

| Statistic | Student ($N$ = 183) | Item ($N$ = 20) |
|---|---|---|
| Separation statistics | | |
| Reliability of separation | 0.84 | 0.94 |
| $\chi^2$ | 234.5* | 135.4* |
| $df$ | 51 | 19 |

*$p < .05$

### Critical Thinking based on Aspect Interpretation

The interpretation aspect includes sub-skills such as categorization, deciphering significance, and clarifying meaning. The results of data analysis using the Rasch Model suggest that students' ability to interpret the concept of work and energy varied depending on their level of understanding (see Figure 1). The analysis indicated that students with a score of 1 had a logit of -3.32, indicating significant difficulty in interpreting the concept. Students with a score of 2 with a logit of -1.88, also demonstrated a lack of understanding, though they performed better than those with a score of 1. However, students with a score of 3 and a logit of 1.06 exhibited stronger interpretation skills. Notably, no students achieved a score of 4 (see Figure 1: Item Characteristic Curve, ICC).

A large proportion of students received scores of 1 and 2. These results indicate that students tended to misinterpret how potential and kinetic energy interact during the pole vaulting process. However, this misunderstanding may reflect weaknesses in the teaching model, which does not sufficiently emphasize the application of physics concepts to real-life situations. An example of a student's response is shown in Figure 1.
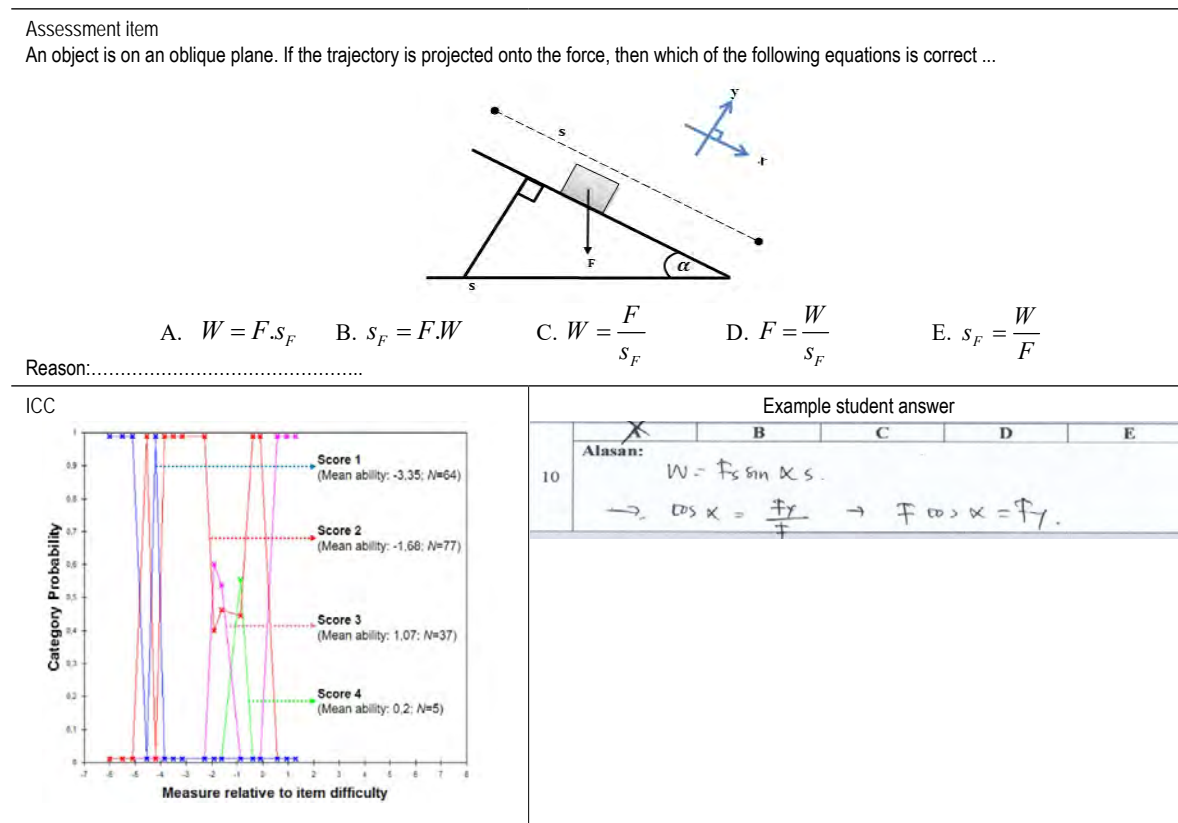
**Figure 1**
Critical Thinking based on Aspect Interpretation



| Assessment item | Item Characteristic Curve (ICC) |
|---|---|
| Consider the following figure! <br><br> The picture is illustrated with an athlete doing a long jump. The statement about mechanical energy is correct ... <br><br> A) Ep at the beginning and end are the same <br> B) Ek at the start and end points are the same <br> C) Ek at the peak is maximal <br> D) Absence of kinetic energy at all points <br> E) Ep at the peak is maximal <br> Reason: ................................................. | Score 1 (Mean ability: -3,32; N=61) <br> Score 2 (Mean ability: -1,88; N=78) <br> Score 3 (Mean ability: 1,06; N=44) |

Example student answer

| | A | B | C✗ | D | E |
|---|---|---|---|---|---|
| 17 | Alasan: Besarnya Energi kinetik di Puncakbennai Maksimal | | | | |

Reason in English (EN):
The magnitude of the kinetic energy at the vertex in the maximal

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

### Critical Thinking based on Aspect Analysis

The analysis aspect refers to the ability to identify intended and actual inferential relationships between statements, questions, concepts, descriptions, or other forms of representation. The analysis aspect was measured by how well students identified and understood the relationships between physics variables, such as the effort exerted by a ball on a curved trajectory. The results of the Rasch model analysis indicated differences in students' abilities in this area (see Figure 2). Students with a score of 1, with a logit of -3.35, showed significant difficulty in analyzing the relationships between physics variables. Students with a score of 2, with a logit of -1.68, demonstrated slightly better analytical skills but still struggled to fully understand and connect the concepts needed to solve the problem. However, students with a score of 3, with a logit of 1.07, exhibited better abilities in correctly analyzing the correlation between physics variables. Students with a score of 4, with a logit of 0.2, showed strong analytical skills and consistently understood and solved problems based on the concepts they had learned (see Figure 2: Item Characteristic Curve, ICC).

**Figure 2**
*Critical Thinking based on Aspect Analysis*



A large number of students received scores of 1 and 2, indicating that their analytical skills were relatively weak. It means that students often had difficulty connecting concepts, which may be attributed to factors such as a lack of foundational understanding of work and energy or ineffective teaching methods in explaining the practical application of these concepts. Few students (N=5) achieved the maximum score of 4, indicating the ability to identify the association between relevant physics variables, such as the relationship between force, distance, and energy change, in the context of a skater moving along a curved track. An example of a student's response is presented in Figure 2.

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

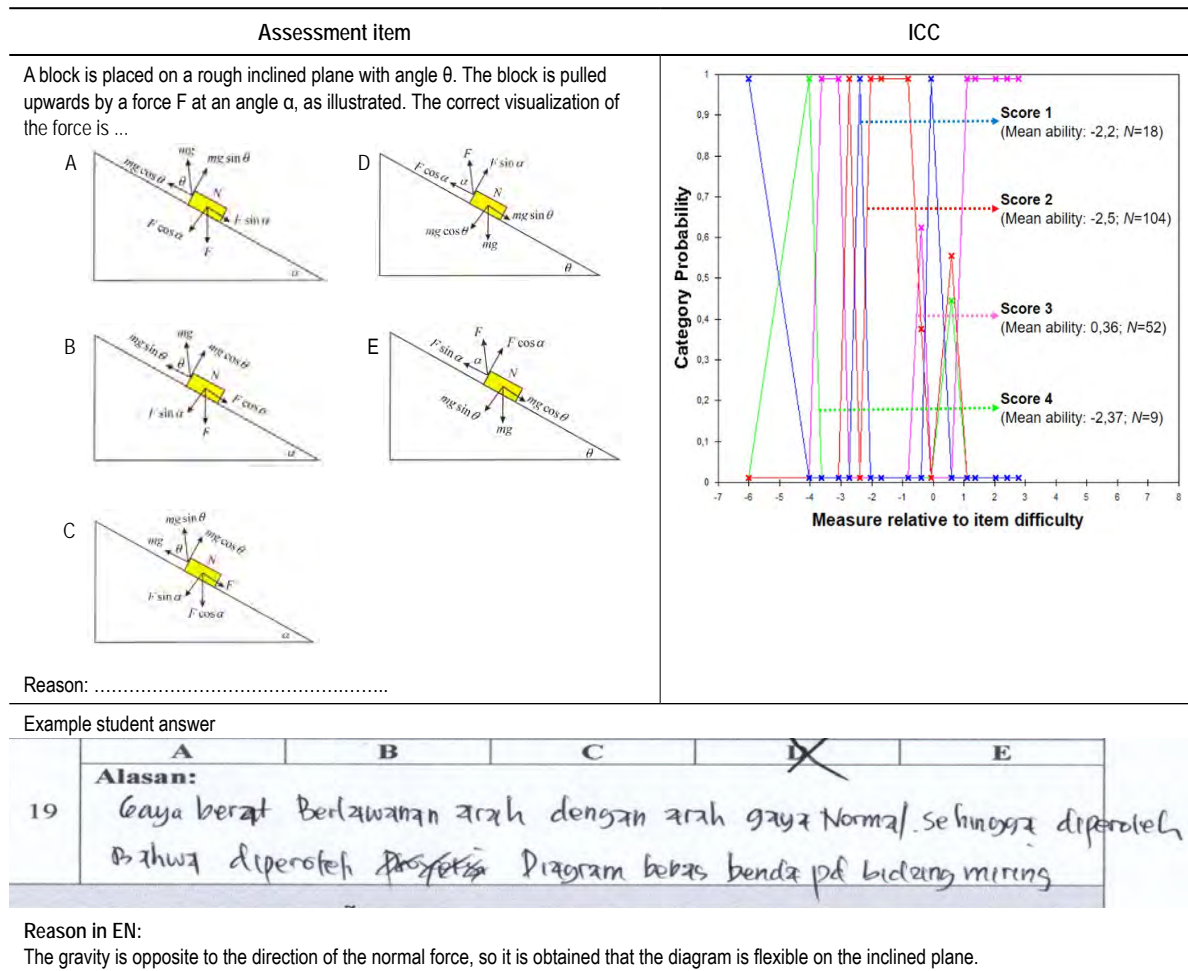ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

### Critical Thinking based on Aspect Evaluation

The evaluation aspect involves assessing the credibility of statements or representations and the logical strength of inferential relationships between them. The results of the Rasch model analysis for the evaluation aspect revealed variations in students' response patterns when distinguishing between valid and invalid information (see Figure 3). For instance, in the context of upward force at a specific angle, the logit analysis indicated that scores of 1 (-2.2 logits), 2 (-2.5 logits), 3 (.36 logits), and 4 (-2.37 logits) reflected differing levels of difficulty in assessing the validity of information (see Figure 3: Item Characteristic Curve, ICC). Lower logit scores indicated that students struggled with evaluating information accurately, while higher scores suggested a stronger ability to evaluate and distinguish between valid and invalid data.

A significant portion of students received scores of 1 and 2, indicating difficulty in projecting the correct direction of force in scenarios where a beam exerts upward force at an angle. The analysis revealed that many students had trouble distinguishing between force components parallel to the surface and those perpendicular to it. Such difficulties highlight the need to strengthen students' understanding of vectors and force projection in physics. One student's response is provided in Figure 3.

**Figure 3**
*Critical Thinking based on Aspect Evaluation*

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

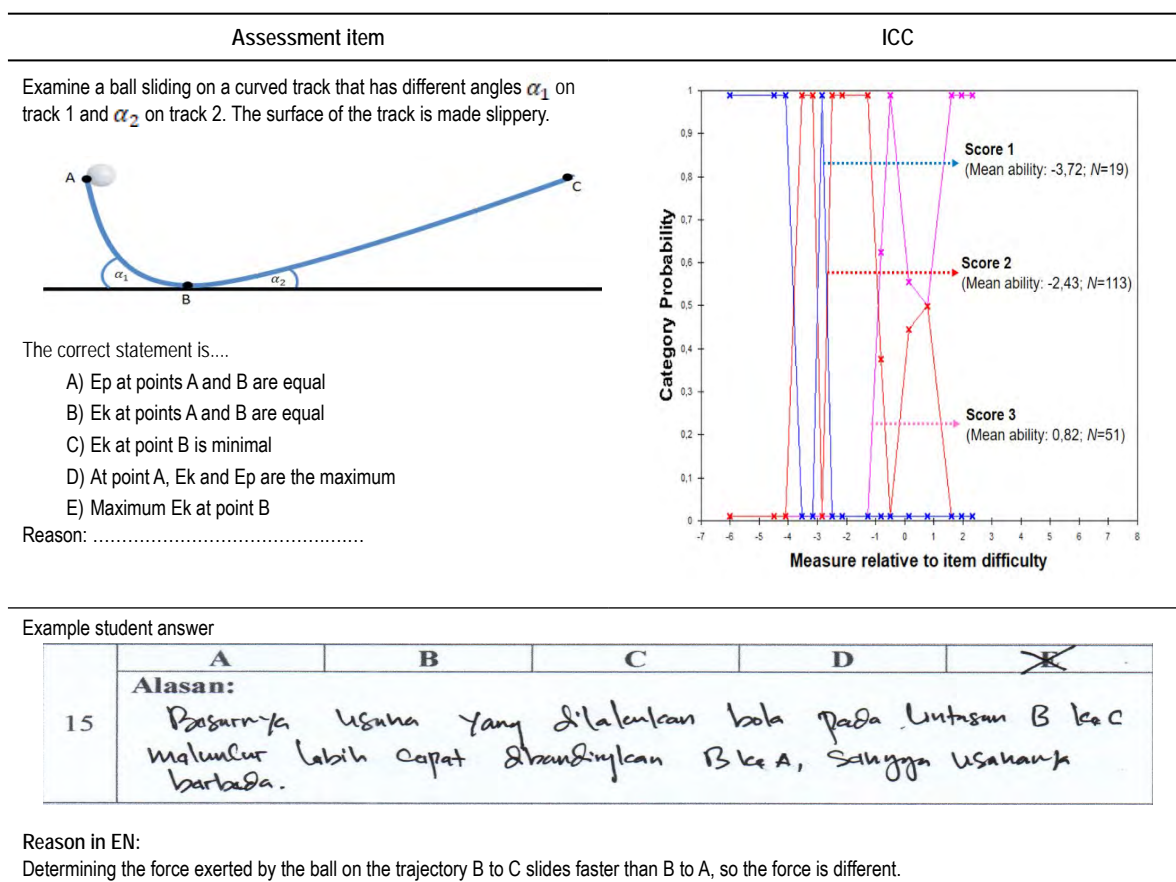ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

### Critical Thinking based on Aspect Explanation

The results of the Rasch model analysis indicate significant variations in students' ability to explain the concepts of work and energy (see Figure 4). Students with a score of 1 and a logit of -3.72 presented great difficulty explaining this concept logically and coherently. Students with scores of 2 and a logit of -2.43 also encountered challenges in explaining the idea despite having a better understanding than students with scores of 1. Meanwhile, students with a score of 3 and a logit of .82 showed a better ability to explain the concept of work and energy. However, no students scored 4 on the explanation aspect (see Figure 4: ItemCharacteristic Curve, ICC).

Most students received scores of 1 or 2, suggesting a widespread struggle to provide complete and precise explanations. The question item related to the explanation aspect, which depicted a ball sliding on a curved track at varying angles, posed a challenge for many students. They had difficulty explaining how changes in the track's angle influenced the work done on the ball and the associated changes in kinetic and potential energy. An example of a student's response is provided in Figure 4.

**Figure 4**
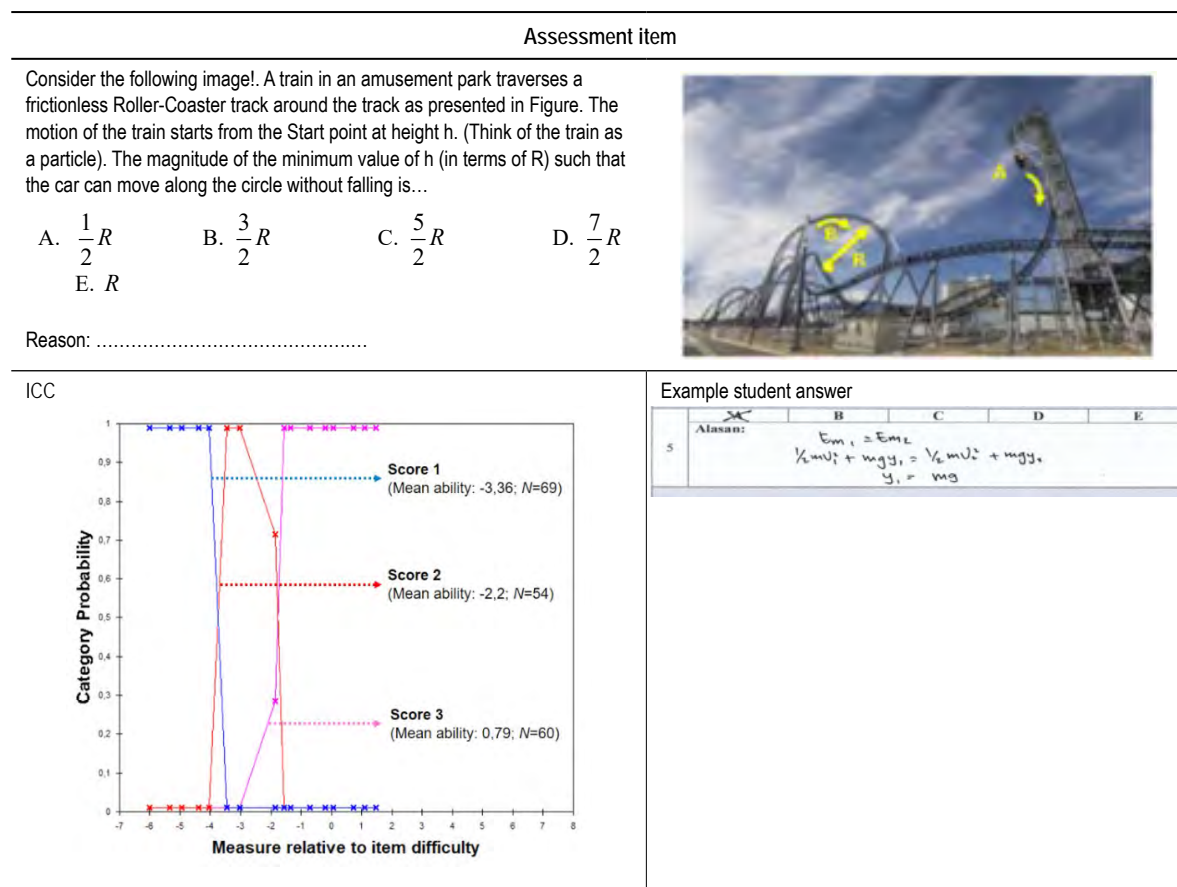*Critical Thinking based on Aspect Explanation*



### Critical Thinking based on Aspect Inference

Aspect inference includes identifying and preserving the elements necessary to make reasonable inferences, forming hypotheses, and considering relevant information to produce logical consequences from existing data, statements, principles, or evidence. The results of the Rasch model analysis indicate variations in students' ability to make appropriate inferences (see Figure 5). Students with a score of 1 and a logit of -3.36 indicated significant difficulty drawing accurate inferences from the given data. Students with a score of 2 and a logit of -2.2 represented

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(PP. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

a slight enhancement in inference ability but still had difficulty connecting the various relevant elements to produce a reasonable conclusion. Furthermore, students with a score of 3 and a logit of .79 displayed a better ability to make inferences. However, no students scored 4 (see Figure 5: Item Characteristic Curve, ICC).

**Figure 5**
*Critical Thinking based on Aspect Inference*



The research findings on the explanation aspect were that most students obtained a score of 1 and a score of 2. It indicates that students have difficulty connecting data with relevant physics concepts, so efforts are needed to improve students' inference skills, focusing on developing data analysis skills and connecting physics concepts through more structured exercises and intensive guidance. The question item on the explanation aspect is presented with a figure of a roller coaster track, where students are asked to plan the minimum height so that the car can travel the circular track ideally. One example of a student's answer is presented in Figure 5.

*Critical Thinking based on Aspect Self Regulation*

An aspect of self-regulation includes a person's ability to monitor their cognitive activities consciously, the elements used in them, and the results obtained. The results of data analysis using the Rasch model demonstrate that students' ability levels vary depending on their level of understanding (see Figure 6). According to the results of the analysis, students with Score 1 of -3.36 logits indicated the lowest self-regulation ability, likely due to the inability to monitor and adjust the approach used in solving the problem. Meanwhile, students with a Score of 2 of -1.54 logit and a Score of 3 of .88 logit indicate a gradually increasing self-regulation ability, where students begin to be able to monitor and evaluate their approach. However, there may still be deficiencies in the self-correction process. Students with a Score of 4 of 2.38 logits indicated the highest self-regulation ability, reflected in the ability to effectively monitor, evaluate, and adjust their approach to solving problems and make necessary corrections to achieve better results (see Figure 6: Item Characteristic Curve ICC).
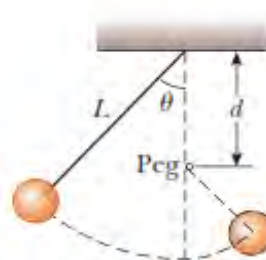
The main findings from measuring self-regulation aspects in the context of understanding and applying the concepts of work and energy indicated a significant correlation between these aspects and students' overall performance. The measurement results reveal that students with better self-regulation skills tend to show a deeper understanding and more precise application of complex physics concepts, such as in the problem of the minimum distance of nails on the trajectory of a mathematical swing. One example of a student's answer is presented in Figure 6.

**Figure 6**
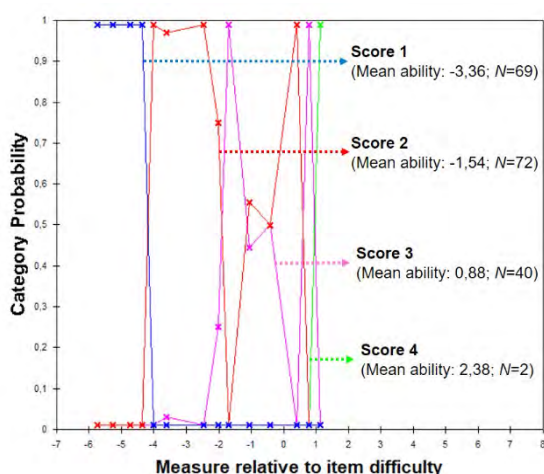*Critical Thinking based on Aspect Self Regulation*



Assessment item

Consider the following image!. A pendulum consists of a string of length L and a small ball swinging in the vertical plane illustrated in the figure. The string touches a peg located at a distance d below the hanging point. If the pendulum is released from rest in a horizontal position $\left(\theta = 90^0\right)$ and will swing in a complete circle centered on the stake, then the minimum value of d is …

A. $\frac{1}{2}L$   B. $\frac{3}{4}L$   C. $\frac{2L}{5}$   D. $\frac{3L}{5}$   E. $L$

Reason: ……………………….............................................

ICC

Score 1 (Mean ability: -3,36; N=69)
Score 2 (Mean ability: -1,54; N=72)
Score 3 (Mean ability: 0,88; N=40)
Score 4 (Mean ability: 2,38; N=2)

Example student answer

**Discussion**

The current study indicates that the two-tier test to measure critical thinking on physics topics, particularly work and energy, has good measurement quality. According to the statistics obtained, the standard error for test takers was .34 logits with a standard deviation of .03, while for the items, the standard error was .14 logits with a standard deviation of .01. The average infit MNSQ value for participants was .92 ($SD = 0.15$), and for items, it was .94 ($SD = 0.17$), indicating that most items had a good level of fit with the Rasch model. The mean square (MNSQ) outfit values also showed promising results, with .97 ($SD = .021$) for participants and .99 ($SD = 0.19$) for items, indicating that there were no significant deviations from the model. The reliability of separation for participants was .84, and for items, it was .94, demonstrating the test's ability to separate participants' skills and identify items that vary in difficulty. The chi-square results were significant ($p < .05$) for both participants and items, confirming that the model could explain the variability in the data well. Unsurprisingly, this finding is consistent with the results of previous studies, where infit and outfit mean square values close to 1 indicate a good level of fit between the data and the Rasch model (Cvenic et al., 2022; Mešić et al.,

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

2019). Similarly, the infit and outfit mean square averages for participants and items indicated no significant deviations from the model. Moreover, the high reliability of separation for both participants and items has been discussed in research Bond, (2015) and Lu, (2024), which emphasizes that high reliability indicates the instrument's ability to clearly distinguish between different groups of participants in terms of ability and variation in item difficulty levels.

The interpretation aspect in the context of physics education refers to students' ability to understand, interpret, and construct meaning from information, concepts, or data. The results indicated that most students had low interpretation skills, with an average score of 1 at -3.32 logits and a score of 2 at -1.88 logits. However, students did not achieve a score of 4, indicating significant limitations in the interpretation aspect. The study's results differ from research conducted by Ennis (1996), (2018), where students trained with methods that encourage reflective thinking tended to have higher interpretation scores. This difference is likely due to the learning model used in the local context, which emphasizes developing interpretation skills less.

The analysis aspect involves the ability of students to break down information into smaller parts to understand the structure and relationships between these components (Facione, 2020). The results indicated that most students had a low level in the analysis aspect, with an average score of 1 at -3.35, 2 at -1.68, 3 at 1.07, and 4 at .2. The findings showed that although some students demonstrated analytical skills, most were still at a low level. According to Paul & Elder (2014), the analysis aspect depends highly on students' educational background and exposure to problems requiring complex solutions. Another study by Halpern (2019), found that problem-based learning can significantly improve students' analytical skills. The findings suggest that the learning methods used were not optimal in developing the analysis aspect, contrasting with Halpern (2019) research, which demonstrated significant improvement after appropriate intervention.

The evaluation aspect involves assessing the credibility of information sources and the strength of the arguments presented (Facione, 2020). The results revealed that most students had low evaluation skills, with an average score of 1 at -2.2 logits, 2 at -2.5 logits, and 3 at .36 logits. However, there was a decrease in score 4 to -2.37 logits, indicating inconsistencies in students' evaluation skills. Facione (2020), argues that evaluation is the most challenging aspect of critical thinking because it requires integrating various other thinking skills. Ennis, (1996), (2018) also found that the evaluation aspect requires continuous practice and proper contextualization in education. The study's findings differ from previous studies due to the lack of implementation of evaluation-focused teaching strategies, such as structured debates or peer-reviewed assessments, which were reported to enhance students' evaluation skills (Ataizi & Donmez, 2014; Kang et al., 2010).

Explanation involves demonstrating one's understanding clearly and thoroughly and supporting arguments with relevant evidence (Facione, 2020). The results indicated that most students had low explanation skills, with an average score of 1 at -3.72 logits, 2 at -2.43 logits, and 3 at .82 logits, indicating that most students struggled to provide adequate explanations. The absence of students who scored 4 suggests that explanation skills need significant improvement. According to a study by Sithole (2023), explanation ability is strongly related to verbal and written communication skills. Another study by Pursitasari et al. (2020), stated that explanation skills can be enhanced through exercises involving critical essays and oral presentations. The results of the study differ from previous studies, where the explanation aspect was more developed. In contrast, previous research suggested that the explanation aspect can be significantly enhanced through structured interventions, such as integrating explanation tasks into the curriculum.

The inference aspect refers to the ability to draw conclusions from existing evidence and make predictions based on observed patterns (Facione, 2020). The results indicated that most students had low inference skills, with an average score of 1 at -3.36 logits and 2 at -2.2 logits, while only a few students achieved a score of 3 at .79 logits. No students scored 4, indicating that the ability to draw conclusions needs improvement. According to Ennis (1996, 2018), inference is one of the most fundamental critical thinking skills, which can be improved through inquiry-based learning and experimentation. Research by Falloon et al. (2022), found that using case studies in learning can help students develop stronger inference skills. The results suggest that students' inference skills were still low, reflecting the lack of practical pedagogical approaches proposed in previous studies.

The self-regulation aspect includes students' ability to regulate their thought processes, monitoring and reassessing their thoughts and actions (Facione, 2020). Interestingly, in this aspect, the results indicated more

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

positive outcomes than in other aspects. However, most students were still at a low level, with a score of 1, which was -3.36 logits, a score of 2, which was -1.54 logits, and a score of 3, which was .88 logits. Some students achieved a score of 4, 2.38 logits, indicating potential for further development. The results showed higher self-regulation abilities in some students, which is in line with research findings by Zimmerman & Schunk (2004), who emphasize the importance of self-regulation as a predictor of long-term academic success. Other research by Cutler and Zimmerman (2011); and Zimmerman and Schunk (2004), suggests that teaching strategies that encourage self-reflection, such as journaling or self-assessment, could improve students' self-regulation skills.

The current study contributes to understanding critical thinking in the context of physics education and science education more broadly. The study's results deepen insights into physics regarding how students understand concepts such as work and energy through the lens of critical thinking. Measuring aspects such as interpretation, analysis, evaluation, explanation, inference, and self-regulation in the study revealed students' critical thinking distribution. It identified areas where students' understanding still needs to be improved. The contribution of the research results is relevant and crucial in understanding abstract concepts in physics. The results also offer a new approach to measuring critical thinking specific to certain topics in physics, such as work and energy, using a logit scale. The results allow for a more accurate and focused evaluation of students' ability to understand and apply physics concepts and can serve as a basis for developing more effective learning models in teaching physics at higher levels.

### Limitations and Directions for Future Research

The current study used a two-tier test calibrated with the Rasch model, but the calibration is only partially optimal for some contexts or populations. Variations in students' interpretation of questions or differences in educational backgrounds and learning experiences may affect the accuracy of the measurement. However, it may also affect the generalizability validity of the results, mainly if applied to a population different from the sample used in this study. Future research should address this limitation by extending the instrument's calibration to various contexts and more heterogeneous populations. Moreover, it is worth considering the development of a more comprehensive instrument to measure aspects of critical thinking that were missed in this study. Using more diverse methodologies, such as longitudinal analysis or mixed approaches, may also provide deeper insights into developing critical thinking over the long term and in various educational settings.

## Conclusions and Implications

The current study is the utilization of a two-tier test calibrated with Rasch analysis, which provides in-depth and detailed insights into students' critical thinking in the context of physics education. The measurement results successfully identified variations in critical thinking skills among students in six aspects: interpretation, analysis, evaluation, explanation, inference, and self-regulation. The results suggest that most students have critical thinking skills that are still low to moderate, with some aspects, such as interpretation and evaluation, showing the need for further development. The absence of students who reached the highest level in some aspects indicates a significant gap in the mastery of critical thinking skills essential for understanding complex concepts in physics. The research's main strength is the ability of the two-tier test and Rasch analysis to provide a more valid and reliable evaluation of students' critical thinking aspects. The study results allow educators and researchers to measure students' overall level of understanding and diagnose specific areas that require further educational intervention. The research findings contribute to an enhanced sense of how critical thinking can be developed and evaluated more effectively in physics learning, as well as provide avenues for further research in developing more precise assessment instruments in physics and science education.

## Declaration of Interest

The authors declare no competing interest.

# References

Affandy, H., Nugraha, D. A., Pratiwi, S. N., & Cari, C. (2021). Calibration for instrument argumentation skills on the subject of fluid statics using item response theory. *Journal of Physics: Conference Series*, *1842*(1), 1–10. https://doi.org/10.1088/1742-6596/1842/1/012032

Affandy, H., Sunarno, W., Suryana, R., & Harjana. (2024). Integrating creative pedagogy into problem-based learning: The effects on higher order thinking skills in science education. *Thinking Skills and Creativity*, *53*, Article 101575. https://doi.org/10.1016/j.tsc.2024.101575

Ataizi, M., & Donmez, M. (2014). Book review: 21st century skills-learning for life in our times. *Contemporary Educational Technology*, *5*(3), 272–274.

Bao, L., & Koenig, K. (2019). Physics education research for 21 st century learning. *Disciplinary and Interdisciplinary Science Education Research*, *1*(2), 1–12. https://doi.org/s43031-019-0007-8

Bond, T. (2015). *Applying the Rasch model: Fundamental measurement in the human sciences* (3rd Ed). Routledge.

Cari, C., Pratiwi, S. N., Affandy, H., & Nugraha, D. A. (2020). Investigation of undergraduate student concept understanding on Hydrostatic Pressure using two-tier test. *Journal of Physics: Conference Series*, *1511*(1). https://doi.org/10.1088/1742-6596/1511/1/012085

Cascella, C., Giberti, C., & Bolondi, G. (2020). An analysis of differential item functioning on INVALSI tests, designed to explore gender gap in mathematical tasks. *Studies in Educational Evaluation*, *64*, Article 100819. https://doi.org/10.1016/j.stueduc.2019.100819

Cetin-Dindar, A., & Geban, O. (2011). Development of a three-tier test to assess high school students' understanding of acids and bases. *Procedia - Social and Behavioral Sciences*, *15*, 600–604. https://doi.org/10.1016/j.sbspro.2011.03.147

Cutler, T. D., & Zimmerman, J. J. (2011). Ultraviolet irradiation and the mechanisms underlying its inactivation of infectious agents. *Animal Health Research Reviews / Conference of Research Workers in Animal Diseases*, *12*(1), 15–23. https://doi.org/10.1017/S1466252311000016

Cvenic, M. K., Planinic, M., Susac, A., Ivanjek, L., Jelicic, K., & Hopf, M. (2022). Development and validation of the conceptual survey on wave optics. *Physical Review Physics Education Research*, *18*(1), 10103. https://doi.org/10.1103/PhysRevPhysEducRes.18.010103

Dessie, E., Gebeyehu, D., & Eshetu, F. (2024). Motivation, conceptual understanding, and critical thinking as correlates and predictors of metacognition in introductory physics. *Cogent Education*, *11*(1), Article 2290114. https://doi.org/10.1080/2331186X.2023.2290114

Dwyer, C. P., & Walsh, A. (2020). An exploratory quantitative case study of critical thinking development through adult distance learning. *Educational Technology Research and Development*, *68*(1), 17–35. https://doi.org/10.1007/s11423-019-09659-2

Ennis, R. H. (1996). Critical thinking dispositions: Their nature and assessability. *Informal Logic*, *18*, 165–182.

Ennis, R. H. (2018). Critical Thinking Across the Curriculum: A Vision. *Topoi*, *37*(1), 165–184. https://doi.org/10.1007/s11245-016-9401-4

Facione, P. A. (2020). Critical Thinking: What It Is and Why It Counts. In *Insight assessment: Vol. XXVIII* (Issue 1). https://www.law.uh.edu/blakely/advocacy-survey/Critical Thinking Skills.pdf

Falloon, G., Forbes, A., Stevenson, M., Bower, M., & Hatzigianni, M. (2022). STEM in the making? Investigating STEM learning in junior school makerspaces. *Research in Science Education*, *52*(2), 511–537. https://doi.org/10.1007/s11165-020-09949-3

García-Carmona, A. (2023). Scientific Thinking and Critical Thinking in Science Education: Two Distinct but Symbiotically Related Intellectual Processes. *Science and Education*, *32*(5), 1221–1225. https://doi.org/10.1007/s11191-023-00460-5

Halpern, D. F. (2019). *Thought & Knowledge: An Introduction to Critical Thinking* (Fourth Edi). Lawrence Erlbaum Associates.

Irmak, M., Inaltun, H., Ercan-Dursun, J., Yaniş-Kelleci, H., & Yürük, N. (2023). Development and application of a Three-Tier Diagnostic Test to assess pre-service science teachers' understanding on work-power and energy concepts. *International Journal of Science and Mathematics Education*, *21*(1), 159–185. https://doi.org/10.1007/s10763-021-10242-6

Istiyono, E. (2016). The application of GPCM on MMC test as a fair alternative assessment model in physics learning. *Proceeding of 3rd International Conference on Research, Implementation and Education of Mathematics and Science*, *May*, 25–30.

Istiyono, E., Mustakim, S. S., Widihastuti, Suranto, & Mukti, T. S. (2019). Measurement of physics problem-solving skills in female and male students by PhysTeProSS. *Jurnal Pendidikan IPA Indonesia*, *8*(2), 170–176. https://doi.org/10.15294/jpii.v8i2.17640

Kaltakci, D., Eryilmaz, A., & McDermott, L. C. (2016). Identifying pre-service physics teachers' misconceptions and conceptual difficulties about geometrical optics. *European Journal of Physics*, *37*(4), Article 045705.

Kang, L. O., Brian, S., & Ricca, B. (2010). Constructivism in pharmacy school. *Currents in Pharmacy Teaching and Learning*, *2*(2), 126–130. https://doi.org/10.1016/j.cptl.2010.01.005

Kaur, R., Mantri, A., Nagabhushan, P., & Singh, G. (2024). Rasch Computing Analysis of Two Tier Concept Inventory to Assess Engineering Students' Conceptual Knowledge. *SN Computer Science*, *5*(5), 643–656. https://doi.org/10.1007/s42979-024-02955-6

Laliyo, L. A. R., Sumintono, B., & Panigoro, C. (2022). Measuring changes in hydrolysis concept of students taught by inquiry model: Stacking and racking analysis techniques in Rasch model. *Heliyon*, *8*(3), e09126. https://doi.org/10.1016/j.heliyon.2022.e09126

Lu, Y. (2024). Independent predictors of family resilience in patients with ischemic stroke: A cross-sectional survey. *Heliyon*, *10*(3), e25062. https://doi.org/10.1016/j.heliyon.2024.e25062

Lukman, Marsigit, Istiyono, E., Kartowagiran, B., Retnawati, H., Kistoro, H. C. A., & Putranta, H. (2021). Effective teachers' personality in strengthening character education. *International Journal of Evaluation and Research in Education*, *10*(2), 512–521. https://doi.org/10.11591/ijere.v10i2.21629

Mafinejad, M. K., Arabshahi, S. K. S., Monajemi, A., Jalili, M., Soltani, A., & Rasouli, J. (2017). Use of Multi-Response format test in the assessment of medical students' critical thinking ability. *Journal of Clinical and Diagnostic Research*, *11*(9), LC10–LC13. https://doi.org/10.7860/JCDR/2017/24884.10607

Memduhoğlu, H. B., & Keleş, E. (2016). Evaluation of the relation between critical-thinking tendency and problem-solving skills of pre-service teachers. *Journal of Educational Sciences Research*, *6*(2), 75–94. https://doi.org/10.12973/jesr.2016.62.5

Mešić, V., Neumann, K., Aviani, I., Hasović, E., Boone, W. J., Erceg, N., Grubelnik, V., Sušac, A., Glamočić, D. S., Karuza, M., Vidak, A., Alihodžić, A., & Repnik, R. (2019). Measuring students' conceptual understanding of wave optics: A Rasch modeling approach. *Physical Review Physics Education Research*, *15*(1), 1–20. https://doi.org/10.1103/PhysRevPhysEducRes.15.010115

Nyirahabimana, P., Minani, E., Nduwingoma, M., & Kemeza, I. (2023). Multimedia-aided technologies for effective learning of quantum physics at the university level. *Journal of Science Education and Technology*, *32*(5), 686–696. https://doi.org/10.1007/s10956-023-10064-x

Paul, R., & Elder, L. (2014). *Consequential Validity: Using Assessment to Drive Instruction* (White Pape). Foundation for Critical Thinking.

Peng, F. (2023). Evaluating critical thinking of English learners using modern technologies and GTMA. *Soft Computing*, *7*. https://doi.org/10.1007/s00500-023-08127-7

Potvin, P., Skelling-Desmeules, Y., & Sy, O. (2015). Exploring secondary students' conceptions about fire using a two-tier, true/false, easy-to-use diagnostic test. *Journal of Education in Science, Environment and Health*, *1*(2), 63. https://doi.org/10.21891/jeseh.99647

Pursitasari, I. D., Suhardi, E., Putra, A. P., & Rachman, I. (2020). Enhancement of student's critical thinking skill through science context-based inquiry learning. *Jurnal Pendidikan IPA Indonesia*, *9*(1), 97–105. https://doi.org/10.15294/jpii.v9i1.21884

Putica, K. B. (2023). Development and validation of a Four-Tier Test for the assessment of secondary school students' conceptual understanding of amino acids, proteins, and enzymes. *Research in Science Education*, *53*(3), 651–668. https://doi.org/10.1007/s11165-022-10075-5

Rapti, S., & Sapounidis, T. (2024). Critical thinking, communication, collaboration, creativity in kindergarten with educational robotics: A scoping review (2012–2023). *Computers and Education*, *210*(April 2023). https://doi.org/10.1016/j.compedu.2023.104968

Sasson, I., Yehuda, I., & Malkinson, N. (2018). Fostering the skills of critical thinking and question-posing in a project-based learning environment. *Thinking Skills and Creativity*, *29*, 203–212. https://doi.org/10.1016/j.tsc.2018.08.001

Sithole, N. V. (2023). The Efficacy of Microteaching in a Teacher Education Programme During the Lockdown at a University in South Africa. *International Journal of Learning, Teaching and Educational Research*, *22*(2), 76–91. https://doi.org/10.26803/ijlter.22.2.5

ten Dam, G., & Volman, M. (2004). Critical thinking as a citizenship competence: Teaching strategies. *Learning and Instruction*, *14*(4), 359–379.

Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, *10*(2), 159–169. https://doi.org/10.1080/0950069880100204

van Laar, E., van Deursen, A. J. A. M., van Dijk, J. A. G. M., & de Haan, J. (2020). Measuring the levels of 21st-century digital skills among professionals working within the creative industries: A performance-based approach. *Poetics*, *81*(April 2019), 101434. https://doi.org/10.1016/j.poetic.2020.101434

Wan, T. (2023). Investigating student reasoning about measurement uncertainty and ability to draw conclusions from measurement data in inquiry-based university physics labs. *International Journal of Science Education*, *45*(3), 223–243. https://doi.org/10.1080/09500693.2022.2156824

Wang, C.-C., & Ho, H.-C. (2024). Development of the imagination-creativity process scale in design. *Thinking Skills and Creativity*, *53*(April 2024), 101545. https://doi.org/10.1016/j.tsc.2024.101545

Wang, Y., Xu, Z.-L., Lou, J.-Y., & Chen, K.-D. (2023). Factors influencing the complex problem-solving skills in reflective learning: results from partial least square structural equation modeling and fuzzy set qualitative comparative analysis. *BMC Medical Education*, *23*(1). https://doi.org/10.1186/s12909-023-04326-w

Wang, Y., & Zhang, X. (2024). A study of the effect of peer assessment on children's critical thinking in a kindergarten craft course. *International Journal of Technology and Design Education*, *34*(4), 1275–1303. https://doi.org/10.1007/s10798-024-09914-5

Journal of Baltic Science Education, Vol. 23, No. 6, 2024

APPLICATION OF RASCH MODEL IN TWO-TIER TEST FOR ASSESSING CRITICAL THINKING
IN PHYSICS EDUCATION
(pp. 1227–1242)

ISSN 1648–3898 /Print/
ISSN 2538–7138 /Online/

Zakwandi, R., Istiyono, E., & Dwandaru, W. S. B. (2024). A two-tier computerized adaptive test to measure student computational thinking skills. *Education and Information Technologies*, *29*(7), 8579–8608. https://doi.org/10.1007/s10639-023-12093-w

Zimmerman, B. J., & Schunk, D. H. (2004). Self regulating intellectual processes and outcomes: A social cognitive perspective. In D. Y. Dai & R. J. Sternberg (Eds.), *Motivation, emotion, and cognition: Integrative perspective on intellectual functioning and development* (pp. 523–549). Erlbaum Associate Publishers.

**Sujiyani Kassiavera**

PhD Student, Doctorate Program of Natural Science Education, Sebelas Maret University (Universitas Sebelas Maret), Jl. Ir. Sutami No.36, Jebres, Kec. Jebres, Surakarta, Central Java, 57126 Indonesia.
E-mail: kassiavera.sujiyani@student.uns.ac.id
ORCID: https://orcid.org/0009-0002-7292-1952

**A. Suparmi**
*(Corresponding author)*

Professor, Lecturer, Physics Department, Sebelas Maret University (Universitas Sebelas Maret), Jl. Ir. Sutami No.36, Jebres, Kec. Jebres, Surakarta, Central Java, 57126 Indonesia.
E-mail: soeparmi@staff.uns.ac.id
Website: https://iris1103.uns.ac.id/profil-0015095205.asm
ORCID: https://orcid.org/0000-0001-8395-4309

**C. Cari**

Professor, Lecturer, Physics Department, Sebelas Maret University (Universitas Sebelas Maret), Jl. Ir. Sutami No.36, Jebres, Kec. Jebres, Surakarta, Central Java, 57126 Indonesia.
E-mail: cari@staff.uns.ac.id
Website: https://iris1103.uns.ac.id/profil-0006036104.asm
ORCID: https://orcid.org/0000-0002-5717-1156

**Sukarmin Sukarmin**

Professor, Lecturer, Department of Physics Education, Sebelas Maret University (Universitas Sebelas Maret), Jl. Ir. Sutami No.36, Jebres, Kec. Jebres, Surakarta, Central Java, 57126 Indonesia.
E-mail: sukarmin67@staff.uns.ac.id
Website: https://iris1103.uns.ac.id/profil-0002086703.asm
ORCID: https://orcid.org/0000-0002-3767-0660