



https://doi.org/10.31331/medivesveteran.v9i1.3605

# Development of Extended Multiple Choice Numeracy Questions Based on The Reeves Model

\*Tasya Faricha Amelia<sup>1</sup>, Effendi Nawawi<sup>2</sup>, Ratu Ilma Indri Putri<sup>3</sup>, Evy Ratna Kartika Waty<sup>4</sup> <sup>1, 2, 3, 4</sup> Sriwijaya University \*Email: <u>tasyafa2001@gmail.com</u>

Received: November 2024. Accepted: December 2024. Published: January 2025.

#### ABSTRACT

This study aims to develop numeracy questions based on the Reeves development model that produces high-quality numeracy questions. The questions are designed in an expanded multiple-choice format, which allows students not only to choose answers but also to provide in-depth logical reasons. This study uses a design research with a systematic approach involving four main stages in the Reeves model: problem identification and analysis, initial design, prototype development, and implementation and reflection. The research methods include expert validation, small group trials, and field trials involving 26 sixth grade students as research subjects. The results of the study showed that 15 numeracy questions were developed, leaving 10 numeracy questions that were of high quality and had high validity, reliability, varying levels of difficulty, and good discriminatory power. The other 5 numeracy questions were not used because they did not meet the criteria for quality items. The implementation of the questions showed that the expanded multiple-choice format was able to provide an in-depth evaluation of students' abilities and help identify weaknesses in students' understanding in more detail. With these results, this study concludes that the Reeves development model is effective for developing high-quality numeracy questions.

**Keywords**: Question Development, Numeracy, Expanded Multiple Choice, Reeves Model

**How to Cite**: Amelia, T., Nawawi, E., Putri, R., & Waty, E. (2025). Development of Extended Multiple Choice Numeracy Questions Based on the Reeves Model. *Journal Of Medives : Journal Of Mathematics Education IKIP Veteran Semarang*, 9(1),

### INTRODUCTION

Numeracy skills are an important element in mathematical literacy that every individual, especially students at the elementary school level, must have. Numeracy skills not only include an understanding of numbers and basic mathematical operations, but also the ability to apply these concepts in reallife contexts (Damayanti et al. 2024; Itu, Soro, and Wewe 2024; Sidiq, Ayudia, and Sarjani 2023). In the modern era, numeracy skills play a key role in daily decision making, both in the world of work, education and society (Dilla Nurfadillah et al. 2024; Fajriyah 2022). Therefore, mastering numeracy from an early age is an indicator of educational success.

However, various studies and reports, such as PISA (Programme for International Student Assessment), show that Indonesian students' numeracy literacy skills are still relatively low compared to other countries (D, Khasanah, and Putri 2022; Putrawangsa and Hasanah 2022: Umaya, Subali, and Januarsi 2024). This problem is evident in fraction material, especially in the multiplication and division operations of fractions with natural numbers. Students often have understanding difficulty in basic concepts, performing calculations, and connecting this material to the context of everyday life (Fajar Rizqi et al. 2023; Febrivanti and Nurjaman 2023: Widyatma and Ramadhani 2024). For example, many students are unable to answer story-based questions, such as counting ingredients for cooking or dividing items by a certain amount, which involve fraction operations.

One of the causes of low student numeracy skills is the lack of quality questions, both in terms of validity, reliability, level of difficulty, and differentiating power (Hayati and Nindiasari 2024; Suariantini, Werang, and Astawan 2023). The questions that are often used only focus on memorization and the final answer without exploring the reasons behind the answer (Yunarti and Amanda 2022). This condition shows the need for innovation in developing numeracy questions in order to evaluate students' abilities in more depth.

One innovation that can be done is the use of an expanded multiplechoice format. This format is different from conventional multiple-choice choice because each answer is accompanied by a logical reason. This requires students not only to choose the correct answer but also to understand and explain their choices. With this format, the evaluation can cover the cognitive aspects of students more comprehensively, while also helping to identify the difficulties faced by students in understanding the concept of numeracy.

To produce quality questions, this study adopted the Reeves development development model, a systematic that approach involves problem identification and analysis, initial design, prototype development and implementation and reflection (R, Apriliya, and Kosasih 2018). With this model, it is expected that research can produce valid, reliable numeracy questions, have a proportional level of difficulty, and adequate discriminating power. The ultimate goal is to provide a quality evaluation instrument to support numeracy learning in schools.

#### METHOD

This study uses a development design or resource design with the Reeves development model. This model involves four systematic stages, namely problem identification and analysis, initial design, prototype development

and implementation and reflection (R et al. 2018). This method is applied to produce numeracy questions in the form of expanded multiple choices that meet the criteria of validity, reliability, level of difficulty and good discriminating power. Data was collected through interviews and test results.

The first stage of the Reeves development model problem is identification and analysis. At this stage, the researcher conducted a needs analysis and literature review so that problems were found regarding numeracy questions that were of poor quality and unable to evaluate students' abilities in depth. So that a solution was obtained to develop expanded multiplechoice questions.

The second stage of Reeves' development model is the initial design. The questions are designed based on contextual discourse text. Each question includes a main question with multiple choice answers and a column for student reasons. So that the assessment is carried out with two approaches, namely in terms of the multiple choice answers and in terms of the answers to the reasons. Multiple choice answers are analyzed using a multiple choice test analysis form where correct answers are worth 1 and wrong answers are worth 0. The answer to the reason is analyzed using a descriptive test analysis form with a 0-5 assessment scale. The following is the reason assessment scale:

Table 1. Scoring Rubric for Reason Answers

Score	Score Rubric
0	There is no answer or a completely wrong answer.
1	There is little effort without using the concepts of multiplication and
	division of fractions by natural numbers correctly.
2	An attempt to answer is visible, but the concept of multiplying and
	dividing fractions by natural numbers used is significantly incorrect.
3	Understanding of the concept of multiplying and dividing fractions by
	natural numbers is evident, but there are some calculation errors.
4	The steps are correct, but there is a small error in the final presentation.
5	All steps and answers are correct.

The third stage of Reeves' development model is prototype development. The items are validated by experts or an expert review is conducted. Validation conducted in an expert review involves assessing content, constructs, and language (Asmara and Sari 2021). After being validated by expert review, the questions are revised and then tested and interviews are conducted on several students or small groups of students to understanding assess the of the questions, the difficulty of the questions, the suitability of the discourse text with everyday life, the

suitability of the images with the discourse text and suggestions.

The fourth stage of Reeves' development model is implementation and reflection. The revised test items based on the results of trials and interviews with a group of students or small groups, were then tested on a larger group of students, namely 26 grade 6 students. The results of the trial that produced two answers, namely multiple choice answers and reason answers, were each analyzed to be assessed in terms of validity, reliability, level of difficulty and discriminatory power.

Validity is calculated using item correlation analysis. The test items are said to be valid if the r table value > r calculated value. The r table used is 3.882. Reliability shows the consistency of the measurement results calculated using the Cornbach Alpha coefficient. The following are the reliability assessment criteria according to Safi'i, Alfi, and Fatih (2024):

Reliability Value	Information
< 0.20	Very Low
0.20-0.40	Low
0.41-0.60	Currently
0.61-0.80	Tall
0.81-1.00	Very high

Table 2. Reliability Criteria

The level of difficulty is indicated by the proportion of students who answer correctly. The following formulas and assessment criteria are used to measure the level of difficulty in multiple-choice answers according to Utami and Adilla (2022):

$$P = \frac{B}{JS}$$

Information : P = Difficulty level B = The number of students whoanswered the questions correctly IS = Total number of students

# Table 3. Multiple Choice AnswerDifficulty Level Criteria

MarkP	Information
< 0.30	Difficult
0.31-0.70	Currently
0.71-1.00	Easy

Meanwhile, the formula and assessment criteria used to measure the level of difficulty of the analysis answers according to Astuti, Waluya, and Asikin (2020):  $TK = \frac{Skor Rata - Rata Suatu Soal}{Skor Maksimum Suatu Soal}$ Information : TK = Difficulty level

#### Table 4. Criteria Level of Difficulty Answer Reason

Mark <i>TK</i>	Information
< 0.30	Difficult
0.31-0.70	Currently
0.71-1.00	Easy

Discriminatory power measures the extent to which a question can differentiate between students with high and low abilities. The assessment between multiple-choice answers and reason answers to calculate their discriminatory power is the same. The following is the formula and criteria for assessing discriminatory power in multiple-choice answers and analysis answers according to Rahmasari and Ismiyati (2016):

> $DP = \frac{\bar{X}KA - \bar{X}KB}{Skor Maksimal}$ Information : DP = Distinguishing Power $\bar{X}KA = \text{Upper group average}$  $\bar{X}KB = \text{Lower group average}$

#### **Table 5. Distinguishing Power Criteria**

MarkP	Information
< 0.19	Not Good/Discarded
0.20-0.29	Enough
0.30-0.39	Good
≥ 0,40	Very good

After being measured in terms of validity, reliability, level of difficulty and discriminating power, it will be determined which questions are of good quality. The criteria for good quality questions are questions that have

content validity according to the curriculum or the questions must be valid, have high reliability (>0.6), have varying levels of difficulty including easy, medium and difficult categories, and adequate discriminating power (>0.2) or the questions have minimal sufficient discriminating power (Sa'idah, Yulistianti, and Megawati 2019).

This research is expected to produce numeracy questions that are not only able to measure students' abilities accurately but also provide a deeper picture of students' way of thinking, especially in understanding and applying numeracy concepts in real life.

## **RESULTS AND DISCUSSION**

The development of multiplechoice numeracy questions was expanded using the Reeves development model consisting of four main stages:

## **Problem Identification and Analysis**

At this stage, the researcher identified and analyzed the problem. Based on the identification and analysis of the problem, it was found that the problem was related to the low-quality numeracy questions and were unable to evaluate students' abilities in depth. So that a solution was obtained to develop expanded multiple-choice questions. The first thing that was done was a needs analysis, the researcher identified the purpose of developing numeracy questions, namely to produce quality numeracy questions to measure students' numeracy abilities validly, reliably, and have a good level of difficulty and differentiating power.

After determining the objectives, the researcher began to conduct a literature review by collecting references on numeracy standards,

multiple-choice expanded question formats and quality question criteria. The researcher also reviewed the concept of validity in terms of construct content and language and also reviewed the concept of reliability, level of difficulty and discriminatory power. The researcher also determined that the indicators of quality numeracy questions to be developed were to have good question criteria in terms of: content validity according to the curriculum, high reliability of more than 0.6, varying levels of difficulty of easy, moderate, difficult and adequate discriminatory power of more than 0.2.

## **Initial Design**

The first thing to do at the design stage is to compile questions. Questions are compiled by determining the competencies being measured such as CP, TP and ATP. Then, the researcher also compiles a question grid containing question indicators and cognitive levels.

After determining the competencies being measured and compiling the question grid. The researcher began to create questions. The compilation of questions was carried out using an expanded multiplechoice format where the answer choices consisted of a combination of options that tested in-depth understanding. Each question included a main question with multiple-choice answers and a column for student reasons. So that the assessment was carried out with two approaches, namely in terms of the multiple-choice answers and in terms of the answers to the reasons. The context used in the questions must also be relevant to the real situation.

The questions consist of 15 expanded multiple-choice questions that have 5 numeracy discourse texts. 1 discourse consists of 3 questions with

different cognitive levels. The following are the questions developed at the initial

design stage:



#### **Prototype Development**

At the prototype development stage, this stage focuses on testing questions and improvements based on the results of empirical analysis (Apriatni, Yuhana, and Sukirwan 2022). After conducting an independent evaluation, the researcher began to consult or validate with an expert review to validate the questions in terms of content, construction, and language. The validators consisted of 3 people who were mathematics teachers who had a master's degree in mathematics education. The following are comments

from the expert review validation and revision:

Validators	Comment	<b>Revisions and Improvements</b>
Validator 1	Questions need to be supplemented with related discourse text; number 1 is better to go straight to the point; numbers 4 and 10 are too vague.	Add a caption to the discourse text; clarify questions 4 and 10 with wording such as: "How many friends got watermelon/pizza?"
Validator 2	Discourse text 1 needs to be aligned; use consistent terms in number 1; use fraction format in multiple choices; number 8 needs to be clarified.	Refine the wording of discourse text 1; use the consistent term "wheat flour"; use fractions, not decimals, in fraction problems; clarify the logic in problem number 8.
Validator 3	Pay attention to the consistency of the terms in number 1; number 6 is ambiguous; the text of discourse 3 and the image do not match; questions 9 and 12 need to be clarified.	Align the terms "wheat flour" and "flour"; correct number 6 regarding the remaining watermelon pieces; adjust the text of the discourse and the picture; clarify questions number 9 and 12.
Afte	er conducting an expert	interview questions includ discussions about understanding t

#### **Table 6. Expert Review Validation Results**

After conducting an expert review, the researcher conducted a trial on a group of students or involved a small sample of students to work on the Prototype questions from the results of the revised export validation review. This small sample or small group consisted of 4 grade 6 students. The trial had been conducted, the researcher also conducted interviews with the four students to ask about the students' responses when working on the questions that were developed. The

ed he questions, the difficulty the of questions, the suitability of the discourse text in the questions with everyday life, the suitability of the images with the discourse text given and suggestions for improvement regarding the questions that were developed. The following are the results of the trial and interviews along with their revisions:

Aspect	<b>Student Comments</b>	<b>Revisions and Improvements</b>
Discourse 1: Making Cakes	The measurement from kg to cups is confusing; the image does not support the text.	Convert measurements to cups; add images that match the text to explain ingredient measurements.
Question Number 1	Not straight to the point, hard to understand.	The wording has been clarified to be more direct and easy to understand, while maintaining the same answer choices.
Question Number 2	It's hard to understand.	The question wording has been revised to make it clearer and more relevant to the discourse; the answer choices have not been changed.
Questions	Question number 7 is	The wording of question number 7 has

#### **Table 7. Trial Results**

Number 7 and 8 unclear, number 8 is considered appropriate		been clarified to make it relevant to the discourse. Number 8 remains intact		
Discourse 5: Making JuiceNot relevant to everyday life; question number 15 is confusing.The amount of juice from the has been clarified; the wording 15 has been adjusted to make understand, the answer choid the same		The amount of juice from the ingredients has been clarified; the wording of number 15 has been adjusted to make it easier to understand, the answer choices remain the same		
Difficulty Level easy		Questions that are too difficult and too easy are revised to increase the ideal level of difficulty ( $P = 0.3 - 0.7$ )		
Distinguishing Power		Questions with poor discriminatory power (<2.0) were corrected or discarded.		

After the revision, the researcher began to enter the next stage regarding implementation and reflection. The following are the final results of the test items after expert review validation and small group trials:

Nama :	Tanggal : A	dasannya :			6. Jika adik Budi mendapat $\frac{3}{4}$ b	agian dari semua potongan semangka, berapa bagian yang didapat adiknya
Kelas :					Bodi? a. 50 baejan	c. 10 bagan
Teks Wacama 1 :					b. 30 bagian	d. 25 bagian
Jawablah pertanyaan 1-3 dengan membaca teks warana 1.					Alasaanya :	
Membuat Kue						
Ibu ingin membuat kue untuk acara ulang tahun. Untuk membuat s	atu adonan kue, ibu biasanya Teks Wi	acana 2 :				
mentruturikan 2 getas takaran tepung tengu. Satu getas takaran ber	itsi – kg tepung tengu. Jawabla	h pertanyaan 4-6 dengan memi	baca teks wacana 2.			
/			Membagi Buah			
🥌 = 📷	1.814			<b>т</b> л	ieks Wacana 3 : awabiah pertanyaan 7-8 dengan e	nembacu toks wuonaa 3. Membeli Kain
<ol> <li>Jika ibu mempunyu 4 kg tepung, berapa banyak adonan yang a 1</li> </ol>	ibu baat? Budi me	emiliki 5 buah semangka. Setia	p semangka dipotong menjadi 8 bagian sama	a besar. Budi ingin		1- 10
$b, \frac{1}{2}$ $d, \frac{1}{2}$	membag	gikan potongan semangka kepa	da teman-temannya.		-	
Alasannya :	4. Jik set	ca settap teman mendapat 2 pot mangka?	ong semangka, berapa banyak teman Bodi yi	ang mendapatkan	<b>1</b>	
	a.	15 c	25		-	
						and the second s
		Alasannya :		b	aju dan rok. Untuk membuat satu	n 2 meter kain nank yang masing-masing akan digunakan untuk meninta i baju dibutuhkan ½ meter kain polos. Untuk membuat satu rok
				d	ibutulakan % meter kain batik.	
<ol> <li>Jika Ibu ingin membual S adonan kue, berapa kg tepung terig a 2 kg.</li> </ol>	u yang dibutuhkan?				<ol> <li>Berapa banyak baju yang bia <ul> <li>A</li> </ul> </li> </ol>	a dibuat dari bahan yang dimiliki ibu sekarang?
b, 6 kg d, 12 kg					b. 5	d. 7
Alasannya :					Alasaanya :	
	5. Jik	ka Budi ingin membuat jus dari	% bagian setiap potongan semangka, berapa	a banyak gelas jus yang		
	bis	sa oronal dari 1 bilah semangka 2	c. 8			
	ь	4	d. 16			
		Masannya :			(max)	
<ol> <li>Ibu memiliki <sup>1</sup>/<sub>4</sub> kg tapang teriga di ramah dan ingin membuat adalah Rp14.000/kg. Berapa unug yang harus ibu bayar uarai</li> </ol>	10 adonan kue. Harga tepung terigu k membeli sisa tepung terigu yang				<ol> <li>Jika harga kain batik Rp25.0 tambahan agar bisa membaa</li> </ol>	00-meter, berapa uang yang harus dibayarkan untuk membeli kain batik 1.16 cok?
kurang? a. Rp126.000 c. Rp216.000	100394 0021 14				a. Rp25.000	c. Rp75.000
b. Rp136.000,- d. Rp96.000,-	L				b. Rp30.000	d. Rp100.000
Alacannya :	11. Jun An kepada a. 2 beg b. 215 b	ich ingin memberikan % bugu 4 temannya, berapa potong p gian sogian	m dari selarp pizza kepinda adıkniye, latis sı itza yang didapat setiap temannya? c. 3 bagian d. 4 bagian	na przemys orbugi rata	<ol> <li>Bandingkan banyaknya jus y banyak?</li> <li>Jeruk, karena menghasilka</li> <li>Apel, karena menghasilka</li> </ol>	ang dihasilkan dari 1 buah jeruk dan 1 buah apel. Jus mana yang lebih n <sup>5</sup> i gelas jus c. Sama banyak 6. Tidak cakup data
					Alasaanya :	
<ol> <li>The Bar manhane Christeler (Coole Denses and and income a</li> </ol>	aten bein anlan dan bein berili senar					
<ol> <li>Ma nu membra o naja dan 12 rok, nenapa persanangaran digunakan ibu?</li> </ol>	12. Jika A	ndi sebelumnya memakan ¾	dari pizza yang dia punya. Kemudian sisar	nya ia bagikan kepada 4		
<ul> <li>a. Baju 2 × lebih banyak</li> <li>c. Sama banyak</li> <li>b. Rok 2 × lebih banyak</li> <li>d. Tidak cukup data</li> </ul>	lemana a 2 bas	iya. Berapa polong pizza yanj zian	g diperoleh temanya Andi? c. 1 hagian	1	5. Jika Budi membuat jus deng	an mencampurkan semua bahan yaitu jeruk dan apel yang Budi punya
Alasannya :	b. <sup>1</sup> / <sub>2</sub> bag	gian	d. 3 bagian		a. 13 gelas	c. 15 gelas
	Alarra	person -	<b>1</b> 2 53		b. 24 gelas	d. 9 gelas
	74850	iniya .			Alasannya :	
Tele Wesser 4						
less wacana 4 :						
awaman pertanyaan 10-12 dengan membaca teks wacama 4.						
Membagi Pizza	Teks Wacan	ui 5 :				
650A 650	Jawablah per	rtanyaan 13-15 dengan memb	aca teles wacana 5.			
	3		Membuat Jus			
	Budi memilil	ici 12 buah jeruk. Satu buah je	rruk menghasilkan % gelas jus. Budi juga i	memiliki 16 buah apel.		
andi membeli 2 buah pizza. Setiap pizza dipotong menjadi 8 bagi	an sama besar. Andi ingin membagikan Satu buah ap	el menghasilkan ½-gelas jus				
izza tersebut kepada teman-temanaya.						
<ul> <li>in ana secap seman mendapat 4 potong pizza, berapa banyak te a. 4 ternari</li> <li>c. 6 ternari</li> </ul>	unan yuen Yang mengahatkan luxya).		yalan 🔵 = Delar			
b. 5 teman d. 7 teman						
Alesanaya :	13. Jika I	Budi ingin membuat 9 gelas j	us jeruk, berapa buah jeruk yang dibutuhka	au?		
	a 10 5 5 5	buah xuah	c. 12 buah d. 4 buah			
	0.00		er i e esta			
	Alasa	nnya :				

#### **Implementation and Reflection**

final At this stage. the researcher's focus is on the application of questions and a definite assessment of whether the questions developed are valid, reliable, have a good level of difficulty and differentiating power or not (Fidia, Puspitawati, and Yakub 2022). The application of the questions was carried out in class 6 at SD Negeri 19 Talang Kelapa which consisted of 26 students. During the application, it was seen that there were students who found it easy to do the numeracy questions and there were students who found it difficult to do the numeracy questions.

After the application or implementation of the questions to students, the results of the application are analyzed to ensure whether the questions are valid and reliable and have a good level of difficulty and differentiating power. Because the test used is an expanded multiple choice where not only multiple choices are assessed and analyzed, but there are logical reasons that must also be assessed and analyzed. So the analysis is carried out using two methods, namely analysis of the form of multiple choice questions and analysis of the answer reasons using the analysis of the essay test form. Both analyzes are carried out by analyzing the validity, reliability of the level of difficulty, and differentiating power. The following are the results of the analysis:

1) Validity

The results of the study on the validity of the questions in the Multiple Choice Analysis showed that 10 questions were declared valid and 5 questions were declared invalid. However, the Reason analysis showed that 15 questions developed were valid. The distribution of 15 questions based on the two validity analyses is as follows:

-	Types	of Analysis	Validity Inde	ex Question Items	Amount	Presentation
-		•	Rcount	2,5,7,9,1	10	66.67%
	Multiple Choice		$\geq 0.338$	0,11,12,13,14,15		
			(valid question	)		
	A	nalysis	Rcount	1,3,4,6,8	5	33.3%
			< 0.338			
			(invalid questi	on)		
			Rcount	1,2,3,4,5,6,7,8,9,10,1	15	100%
	Reaso	n Analysis	$\geq 0.338$	1,12,13,14,15		
-			(valid question	)		
Rce	ount	0	0	0% Multiple		
< 0	.338			Choice	0.66	Tall
(inv	alid			Analysis		
ques	stion)			Reason		
				Analysis	0.88	Very high

#### Table 8. Item Distribution Based on Item Validity

#### 2) Reliability

The following are the results of the reliability values in both analyses:

# Table 9. Item Distribution Based on<br/>Reliability

|--|

3) Difficulty Level

The following is the distribution of 15 questions based on their level of difficulty in both analyses:

# Table 10. Distribution of Items Based onDifficulty Level

Type Diffic Question Amo Present	Туре	Diffic	Question	Amo	Present
----------------------------------	------	--------	----------	-----	---------

s of	ulty	Items	unt	ation
Anal	Index			
ysis				
Multi	Easy	4,6,7,10,13,1	6	40%
ple		4	0	4070
Ĉhoi	Curre	1 2 2 5 11 15	6	40%
ce	ntly	1,2,3,3,11,13	0	
Anal	Diffic	9.0.12	2	20%
ysis	ult	8,9,12	3	
Reas on Anal ysis	Easy	4,5,6,10,14	5	33.3%
	Curre	1,2,3,7,11,12	0	53.3%
	ntly	,13,15	ð	
	Diffic		0	13.3%
	ult	8.9	2	

4) Distinguishing Power

The following are the results of the discriminatory power analysis based on multiple choice and reasons:

Distribution of Essay Question Items Based on Distinguishing Power

#### Table11. Distribution of Grains Based on Discriminating Power

Types of Analysis	Distinguishing Power Index	Question Items	Amount	Presentation
	Not good	3,4,5,6,8	4	26.6%
Multiple Choice	Enough	7,14,15	3	20%
Analysis	Good	1,9,10,12,13	5	33.3%
	Very good	2,5,11	3	20%
	Not good	1,3,8	3	20%
	Enough	2,4,9,14	4	26.6%
Reason Analysis	Cood	6,7,11,12,13,1	6	40%
_	Good	5		
	Very good	5.10	2	13.3%

Based on the four analyses, the following is a combination of validity,

reliability, level of difficulty and differentiating power analyses:

#### **Table 12. Combined Analysis**

	Validity		Reliability		Difficulty Level		Distinguishing Power	
No.	Multipl		Multipl				Multipl	
Questio	e	Reason	e	Reason	Multiple	Daacon	e	Reason
n	Choice	Analysi	Choice	Analysi	Choice	Analysis	Choice	Analysi
	Analysi	S	Analysi	S	Analysis		Analysi	S
	S		S				S	
1	Invalid	Valid		Very	Currentl	Currentl	Good	Not
				high	У	У	0000	good
2 V	Volid	Valid	Tall	Very	Currentl	Currentl	Very	Enough
	v allu			high	У	У	good	Ellough
3	Invalid	Valid		Very	Currentl	Currentl	Not	Not
				high	У	У	good	good
4	Invalid	Valid		Very high Easy	Foor	Not	Enough	
					Lasy	Easy	good	Enough
5	Valid	Valid	Tall	Very	Currentl	Easy	Very	Very

				high	У		good	good
6	Invalid	Valid		Very high	Easy	Easy	Not good	Good
7	Valid	Valid	Tall	Very high	Easy	Currentl y	Enough	Good
8	Invalid	Valid		Very high	Difficult	Difficult	Not good	Not good
9	Valid	Valid	Tall	Very high	Difficult	Difficult	Good	Enough
10	Valid	Valid	Tall	Very high	Easy	Easy	Good	Very good
11	Valid	Valid	Tall	Very high	Currentl y	Currentl y	Very good	Good
12	Valid	Valid	Tall	Very high	Difficult	Currentl y	Good	Good
13	Valid	Valid	Tall	Very high	Easy	Currentl y	Good	Good
14	Valid	Valid	Tall	Very high	Easy	Easy	Enough	Enough
15	Valid	Valid	Tall	Very high	Currentl y	Currentl y	Enough	Good

Quality numeracy questions are questions that have content validity according to the curriculum or the questions must be valid, high reliability (>0.6), varying levels of difficulty including easy, medium, and difficult categories, and adequate discriminatory power (>0.2) or the questions have sufficient minimal discriminatory power. This study shows the importance of ensuring the quality of numeracy questions through analysis of validity, reliability, level of difficulty, and discriminatory power. Based on the combined analysis table, it can be seen that out of 15 questions developed, 10 questions meet the criteria as quality questions based on content validity according to the curriculum, high reliability, distribution of varying levels difficulty, adequate of and discriminatory power. The 10 questions are numbers 2, 5, 7, 9, 10, 11, 12, 13, 14, and 15. The other 5 numeracy questions were not used because they did not meet the criteria for quality items.

Content validity is an important

element because it ensures that each question reflects the competency measured according to the curriculum. In this study, validity was strengthened through multiple choice and reason analysis, which provided a deeper measurement dimension. Of the 15 questions developed, only 10 were considered valid based on the multiple choice analysis. However, in the reason analysis, all 15 questions were declared valid. This shows that the question format that includes reasons can provide additional, more comprehensive an dimension. measurement Content validity is also often reviewed through expert review to ensure its relevance to curriculum standards and the abilities being measured. A study by Damayanti, Daryono, and Hari Rayanto (2023)states that content validity is an aspect that must be evaluated by experts in the relevant field so that questions can be ensured to be relevant to learning objectives.

In addition, high item reliability indicates instrument consistency, so that the measurement results remain accurate even when used repeatedly. In

this study, item reliability reached a high value (0.66 for multiple choice) to very high (0.88 for reasons). This value indicates that the item instrument can be used repeatedly with consistent results under similar conditions. Previous studies, such as those conducted by Kania et al. (2023), revealed that high reliability is essential to ensure the accuracy of evaluation. This is especially true in the context of primary education, where evaluation instruments must be able to capture student performance consistently.

The distribution of difficulty levels is also a concern, with questions spanning easy, medium, and hard categories. This ensures that the questions can effectively evaluate students with different levels of ability. Questions that are too easy or too hard are avoided to maintain balance and relevance of the evaluation. An even distribution ensures that the questions can evaluate students with different levels of ability. For example, questions that are too easy tend not to be able to distinguish students with high ability from those with low ability, while questions that are too hard can hinder student motivation. This is in line with the view Utami and Adilla (2022), which states that an ideal question should have a difficulty index between 0.3 and 0.7 to ensure diversity.

Discriminating power is also a key indicator, where questions with low discriminating power values are improved or removed so that the instrument can differentiate students with diverse understandings. In this study, questions with low discriminating power were improved or removed to maintain the quality of the instrument. Sufficient discriminating power (>0.2) ensures that the evaluation instrument can identify various levels of student ability. According to Susilowati et al. (2024), discriminatory power is a critical aspect in item development because it provides insight into the effectiveness of the item in assessment.

This research is in line with various studies that emphasize the importance of item analysis. For example, Magdalena et al. (2021) showed that validity, reliability, level of difficulty, and discriminatory power are key parameters to ensure that test questions are effective as evaluation tools. This is supported by the results of the analysis which highlighted that 10 out of 15 questions met the criteria as high-quality questions, with adequate levels of validity, reliability, difficulty, and discriminatory power. Another study by Cantika et al. (2024) also found that evaluation instruments that do not meet any of these criteria can reduce the quality of evaluation results. Therefore, it is important to ensure that the questions are not only in accordance with the curriculum but also able to provide meaningful information about student abilities.

## CONCLUSION

study successfully This developed numeracy questions based on the Reeves development model with an expanded multiple-choice format. The results of the study revealed that out of 15 numeracy questions developed, only 10 questions met the criteria for good quality, namely having high validity and reliability, varying levels of difficulty, and adequate discriminatory power. This expanded multiple-choice format allows students not only to choose the correct answer, but also to provide in-depth logical reasons. providing a comprehensive evaluation students' understanding of and numeracy abilities. The Reeves development model has proven

effective in producing valid, reliable, and contextual numeracy questions. This study provides significant contributions for teachers and researchers in designing more quality in-depth numeracy learning and evaluation tools.

## REFERENCE

Apriatni, Sri, Yuyu Yuhana, and Sukirwan Sukirwan. 2022. "Pengembangan Instrumen Literasi Numerasi Materi Trigonometri Kelas X Sma." *EDU-MAT: Jurnal Pendidikan Matematika* 10(2):185. doi:

10.20527/edumat.v10i2.13720.

- Asmara, Adi, and Debby Juita Sari. 2021. "Pengembangan Soal Aritmetika Sosial Berbasis Literasi Matematis Siswa SMP." Jurnal Cendekia: Jurnal Pendidikan Matematika 5(3):2950–61. doi: 10.31004/cendekia.v5i3.982.
- Astuti, Astuti, Stevanus Budi Waluya, and Mohammad Budi Asikin. 2020. "Instrumen Kemampuan Berpikir Kreatif Matematika Untuk Siswa Kelas IV Sekolah Dasar." *Musamus Journal of Primary Education* 3(1):27–34. doi: 10.35724/musjpe.v3i1.3117.
- Cantika, Putu Maha Iswari Putri, Luqman Hakim, Vivi Pratiwi, and Citra Auliyatun Nazwah. 2024. "PELAJARAN **EKONOMI** BISNIS DAN ADMINISTRASI UMUM **KELAS** XI SMK MENGGUNAKAN PROGRAM ANATES VERSI 4 . 0." Jurnal *Multidisipliner* Ilmiah Kajian 8(11):324-41.
- D, Darwanto, Mar'atun Khasanah, and Anggi Monica Putri. 2022. "Penguatan Literasi, Numerasi, Dan Adaptasi Teknologi Pada Pembelajaran Di Sekolah." Eksponen 11(2):25-35. doi:

10.47637/eksponen.v11i2.381.

- Damayanti, Ayu Maya, Daryono Daryono, and Yudi Hari Rayanto. 2023. *Evaluasi Pembelajaran*. Cetakan Pe. edited by E. Tresna Setiawan. CV. Basya Media Utama.
- Zastia Damayanti, Elvina, Dinda Alviony, Aprilia Dwi Putri, Rita Feni, Program Studi, Pendidikan Matematika, Program Studi. Pendidikan Ekonomi, and Program "Sosialisasi Dan Studi. 2024. Pelatihan Metode Jarimatika Untuk Meningkatkan Kemampuan Numerasi Siswa Sekolah Dasar." 4(5):490-95.
- Dilla Nurfadillah, Firyal Nasywa Aufa, and Ichsan Fauzi Rachman. 2024. "Membangun Kualitas Pendidikan Melalui Kemampuan Literasi Dan Numerisasi Dalam Implementasi Kurikulum Merdeka." *ALFIHRIS*: *Jurnal Inspirasi Pendidikan* 2(3):128–40. doi: 10.59246/alfihris.v2i3.876.
- Fajar Rizqi, Ardhian, Bilqis Luthfi Adilla, Erani Sulistiyawati, and Taufiqurrohmah. 2023. "Analisis Kesulitan Belajar Matematika Pada Siswa Sekolah Dasar Dan Alternatif Pemecahannya." Jurnal Pendidikan Dasar Flobamorata 4(1):481–88. doi: 10.51494/jpdf.v4i1.588.
- Fajriyah, Euis. 2022. "Kemampuan Literasi Numerasi Siswa Pada Pembelajaran Matematika Di Abad 21." *Seminar Nasional Pendidikan* 21:403–9.
- Febriyanti, Natasya, and Asep Rudi Nurjaman. 2023. "Analisis Kesalahan Siswa Dalam Menyelesaikan Soal Cerita Matematika Di Sekolah Dasar." *Https://Ejournal.Iaifa.Ac.Id/Index. Php/Dirasah Accepted:* 6(2):322– 28.

- Fidia, Farah, Rinie Pratiwi Puspitawati, Pramita Yakub. 2022. and "Pengembangan Instrumen Soal Higher Order Thinking Skills (HOTs) Materi Jaringan Dan Organ Pada Tumbuhan Kelas XI SMA." Berkala Ilmiah Pendidikan *Biologi* (*BioEdu*) 11(3):745–54. doi: 10.26740/bioedu.v11n3.p745-754.
- Hayati, Rohimatul, and Hepsi Nindiasari. 2024. "Pengembangan Instrumen Asesmen Kompetensi Minimum Numerasi Pada Domain Data Dan Ketidakpastian Untuk Siswa SMP." 6(2):84–93.
- Itu, Maria Alexandria, Viorentina Meo Soro, and Melkior Wewe. 2024. "Profil Kemampuan Numerasi Siswa Sekolah Dasar Di SDK Kisanata." *Polinomial : Jurnal Pendidikan Matematika* 3(2):107– 13. doi: 10.56916/jp.v3i2.921.
- Kania, Nia, Rinovian Rais, Yance Manoppo, Zuli Nuraeni, Ahmad Ahmad, Abdul Munif, Yeti Sulfiati, Anas Anas, Mursidin Mursidin, Reni Suwenti, Geubrina Maghfirah, Hanida Listiani, and Dedi S. 2023. Evaluasi Pendidikan (Sebuah Tinjauan Kritis). CV. Edupedia Publisher.
- Magdalena, Ina, Septy Nurul Fauziah, Nur Faziah. and Siti Fika Sulaehatun Nupus. 2021. "Analisis Reliabilitas, Validitas, Tingkat Kesulitan Dan Daya Beda Butir Soal Ujian Akhir Semester Tema 7 Kelas Iii Sdn Karet 1 Sepatan." BINTANG : Jurnal Pendidikan Dan Sains 3(2):198-214.
- Putrawangsa, Susilahudin, and Uswatun Hasanah. 2022. "Analisis Capaian Siswa Indonesia Pada PISA Dan Urgensi Kurikulum Berorientasi Literasi Dan Numerasi." Jurnal Studi Pendidikan Dan Pembelajaran 1(1):1–12.

- R, Wili Karlina, Seni Apriliya, and E. Kosasih. 2018. "Pengembangan Soal Pada Pembelajaran Penulisan Huruf Kapital Dalam Teks Cerita Pendek." Jurnal Ilmiah Pendidikan Guru Sekolah Dasar 5(4):169–77.
- Rahmasari, Dias, and Ismiyati. 2016. "Analisis Butir Soal Mata Pelajaran Pengantar Administrasi Perkantoran." *Economic Education Analysis Journal* 5(1):317–30.
- Sa'idah, Nusrotus, Hayu Dian Yulistianti, and Eka Megawati. 2019. "Analisis Instrumen Tes Higher Order Thinking Matematika Smp." Jurnal Pendidikan Matematika 13(1):41-54. doi: 10.22342/jpm.13.1.6619.41-54.
- Safi'i, Mochammad Dwi Irvan, Cindya Alfi, and Mohamad Fatih. 2024. "Jurnal Perseda." *Jurnal Persada* VII(2):195–205.
- Sidiq, Fadhil, Inge Ayudia, and Tri Mustika Sarjani. 2023. "Optimalisasi Gerakan Literasi Sekolah Melalui Desain Kelas Literasi Numerasi Di Sekolah Dasar Kota Langsa." Journal of Human and Education 3(3):69–75.
- Suariantini, N. N. G., B. R. Werang, and I. G. Astawan. 2023. "Instrumen Asesmen Numerasi Online Menggunakan Aplikasi Kahoot Pada Mata Pelajaran Matematika Kelas IV Sekolah Dasar." Innovative: Journal Of ... 3(2):5712-5824.
- Susilowati, Dwi Andrianik, Nidaria Saputri, Khafidhotun Ni'mah, Luqman Hakim, and Vivi Pratiwi. 2024. "PENERAPAN APLIKASI ANATES UNTUK MENGANALISIS KUALITAS SOAL HOTS PADA FASE-E ELEMEN PENGGUNAAN APLIKASI PENGOLAH ANGKA ( SPREADSHEET ) DALAM

EVALUASI PEMBELAJARAN : STUDI KASUS PADA SISWA KELAS X AKUNTANSI DI SMK NEGERI 1 KENDAL." 8(11):272– 83.

- Umaya, F., B. Subali, and T. D. "Peningkatan Januarsi. 2024. Literasi Numerasi Siswa Materi Usaha, Energi Dan Pesawat Sederhana Melalui Model Pembelajaran Discovery Learning Pada Kelas VIII SMPN 16 ...." ... Nasional Pendidikan Dan ... 649-57.
- Utami, Lisa, and Raysha Adilla. 2022. "ANALISIS **KETERAMPILAN** PROSES SAINS **SISWA MENGGUNAKAN** VIRTUAL LABORATORY PHYSICS **EDUCATION** TECHNOLOGY PADA (PhET) MATERI **INDIKATOR** ASAM BASA." Journal of Research and Education Chemistry 4(1):50. doi: 10.25299/jrec.2022.vol4(1).9348.
- Widyatma, Yollanda Vannesicha, and Amanda Diva Hadi Ramadhani. 2024. "Analisis Kemampuan Pemecahan Masalah Matematis Pada Materi Bilangan Dan Aljabar Siswa Kelas IV SDN 4 Piji." Jurnal Pendidikan Dan Pembelajaran 3(01):335–49.
- Yunarti, Tina, and Ari Amanda. 2022. "Pentingnya Kemampuan Numerasi Bagi Siswa." Seminar Nasional Pembelajaran Matematika, Sains Dan Teknologi 2(1):44–48.