# A Four-Tier Diagnostic Instrument in Acid-Base Properties of Salt Solution: Development Procedure

**Rima Nuraini**

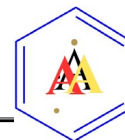Department of Chemistry, Universitas Negeri Malang, Indonesia 65145

*Corresponding author: rimaanuraini@gmail.com*

**Abstract:** This study aims to determine the feasibility of a four-tier diagnostic instrument on salt hydrolysis. The development uses a four-tier diagnostic instrument development procedure by Habiddin & Page (2019) with 6 stages: concept identification, initial test and interview, identification of unscientific student concepts, development of a four-tier diagnostic instrument prototype, prototype validation, and final prototype improvement. The four-tier diagnostic instrument was developed from a multiple-choice instrument open to reasons for capturing student concepts. At last, the finding from this research and development obtained the final product in the form of a four-tier diagnostic instrument with 27 questions that have four levels (tier), the first tier is in the form of questions and answers, second tier is in the level of confidence in the answer chosen, third tier is in the form of selecting the first tier, and the fourth tier is the level of confidence in the reasons chosen. The level of confidence is measured on a scale of 1-5. The instrument developed has an average content validity of 89.45%, with a very decent category and very high reliability (0.858). This shows that the developed four-tier diagnostic instrument is highly feasible for identifying students' misconceptions about salt hydrolysis material.

**Keywords:** four-tier diagnostic instruments, misconceptions, salt hydrolysis, students' understanding

## INTRODUCTION

Salt hydrolysis is one of the chemistry topics taught in 11th-grade high school, according to the 2013 Curriculum. The complex nature of this material lies in the interconnectedness of the concepts being studied with previous concepts. To understand salt hydrolysis well, students are required to understand reaction equilibrium, the dissociation process, and the acid-base properties of reactants and products (Orwat et al., 2017). Additionally, salt hydrolysis is one of the most essential topics in the field of acid-base reactions, yet it is often misunderstood (Secken, 2010). Misconceptions are widely held understandings that do not align with scientific experts' understanding (Pesman & Eryilmaz, 2010). These misconceptions are generally

very difficult to change and can persist for a long time, especially if the teacher-designed classroom learning does not facilitate conceptual change (Demircioğlu et al., 2005). Misconceptions that occur in students during learning can hinder their complete understanding of the material.
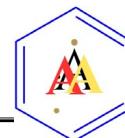
The research results of Maratusholihah et al. (2017) state that 28.12% of students consider salt hydrolysis to be a reaction between water and salt cations or anions, producing $H_3O^+$ and OH- ions, because water breaks down the salt into its cations and anions. Additionally, 18.75% of students believe that salts derived from strong acids and weak bases are acidic because they undergo anion hydrolysis, producing $H_3O^+$ ions, thus increasing the concentration of $H_3O^+$ ions in water. Furthermore, Orwat et al. (2017) reported that 92% of students correctly stated that $ZnCl_2$ solutions are acidic, but the reaction equations they wrote were incorrect. Based on his research, 55% of students noted that the $MgCl_2$ solution was neutral, and 38% correctly stated that $MgCl_2$ was acidic. Most students who answered correctly wrote the hydrolysis reaction with $Mg(OH)_2$ as a product, with 70% of them using a one-way arrow ($\rightarrow$). In comparison, 10% of the students who answered correctly wrote the hydrolysis reaction with $MgOH^+$ as a product. Based on the description, it can be concluded that students have not fully mastered the material on salt hydrolysis.

Students' misconceptions can be identified using several methods, including interviews (Osborne & Gilbert, 1980), concept maps (Novak, 1990), open-ended tests (Taber, 1999), multiple-choice tests (Beichner, 1994), short answer (Billah et al., 2024), Multi-tier instrument (Amala & Habiddin, 2022; Ardina & Habiddin, 2023; Gurel et al., 2015, 2017; Habiddin & Page, 2023; Laliyo et al., 2021) and others. Each instrument used to identify these misconceptions has its own advantages and disadvantages. Among the various methods for identifying misconceptions, the four-tier diagnostic instrument is effective. This test component consists of the first level, which is questions and answers with distractors; the second level is the confidence level of the answers at the first level; the third level is the reason for the answers at the first level; and the fourth level is the confidence level for the chosen reason (Gurel et al., 2017). This four-tier diagnostic instrument allows students to express their different levels of confidence in their answers and reasons, so that students' understanding level can be accurately determined (Habiddin & Page, 2019).

Research on misconceptions regarding salt hydrolysis material has been conducted by Orwat et al. (2017) using questions with four competency tasks, Amelia et al. (2014) using the CRI technique, (Tuysuz, 2009; Ulfah et al., 2024) using a two-tier diagnostic instrument, and Seçken (2010) using multiple-choice and open-ended tests. Based on the literature, no prior research has examined the identification of students' misconceptions about salt hydrolysis using a four-tier diagnostic instrument. Given the advantages of the four-tier diagnostic instrument as described, it is hoped that it will be easier to identify students' understanding of salt hydrolysis.

## METHOD

The development of the four-tier diagnostic instrument in this study adapts the procedure developed by Habiddin & Page (2019) based on the two-tier diagnostic instrument development procedure by Treagust (1988), with modifications to suit. There are six stages involved in developing a four-tier diagnostic instrument: (1)

Concept mapping, (2) Initial testing and interviewing, (3) Identifying students' unscientific concepts, (4) Developing a prototype four-tier diagnostic instrument, (5) Validating the prototype, and (6) Refining the final prototype.

The research subjects for the initial test were students from class XI of SMAN 2 Pare, including classes XI IPA 1, XI IPA 3, and XI IPA 5, totalling 96 students. The research subjects for empirical validation were students from class XI of SMAN 2 Pare, including classes XI IPA 6 and XI IPA 7, totalling 71 students. Content validation was carried out by 1 chemistry lecturer and 2 high school chemistry teachers. The instrument used during the initial test was 30 open-ended multiple-choice questions. The instrument used during empirical validation was a 28-question four-tier diagnostic instrument.

The instrument used for content validation of the four-tier diagnostic instrument was a validation questionnaire with ten assessment indicators. Data analysis techniques include content validation, data analysis, and empirical validation. Empirical validation analysis includes test reliability analysis, item difficulty level, item discrimination power, distractor effectiveness, and item validation. An empirical validation analysis was conducted for each tier: A tier (Answer), R tier (Reason), and B tier (Both).

## RESULTS AND DISCUSSION

### Short Answer Question

*Reliability*

The test reliability was 0.863, indicating that the test items are highly reliable and can be used to develop a four-tier diagnostic instrument.

*Validity*

The validity test results show that 27 items are valid and 3 are not, namely items 5, 20, and 24. The invalid items are considered for revision.

**Table 1.** Validity of short answer questions

| No | R | category | No | R | category | No | R | category |
|----|-------|----------|----|-------|----------|----|-------|----------|
| 1 | 0.655 | Valid | 11 | 0.550 | Valid | 21 | 0.351 | Valid |
| 2 | 0.569 | Valid | 12 | 0.440 | Valid | 22 | 0.417 | Valid |
| 3 | 0.488 | Valid | 13 | 0.334 | Valid | 23 | 0.460 | Valid |
| 4 | 0.457 | Valid | 14 | 0.511 | Valid | 24 | 0.001 | Invalid |
| 5 | 0.191 | Invalid | 15 | 0.496 | Valid | 25 | 0.440 | Valid |
| 6 | 0.327 | Valid | 16 | 0.361 | Valid | 26 | 0.515 | Valid |
| 7 | 0.571 | Valid | 17 | 0.414 | Valid | 27 | 0.474 | Valid |
| 8 | 0.558 | Valid | 18 | 0.528 | Valid | 28 | 0.571 | Valid |
| 9 | 0.422 | Valid | 19 | 0.369 | Valid | 29 | 0.543 | Valid |
| 10 | 0.658 | Valid | 20 | 0.070 | Invalid | 30 | 0.663 | Valid |

*Difficulty Level (P)*

Table 2 shows that there are 6 easy questions, 22 moderate questions, and 2 difficult questions.

**Table 2.** Difficulty level of short answer questions

| No | P | Category | No | P | Category | No | P | Category |
|----|------|----------|----|------|----------|----|------|----------|
| 1 | 0.58 | moderate | 11 | 0.30 | difficult | 21 | 0.52 | moderate |
| 2 | 0.71 | easy | 12 | 0.46 | moderate | 22 | 0.71 | easy |
| 3 | 0.41 | moderate | 13 | 0.88 | easy | 23 | 0.64 | moderate |
| 4 | 0.48 | moderate | 14 | 0.55 | moderate | 24 | 0.39 | moderate |
| 5 | 0.95 | easy | 15 | 0.52 | moderate | 25 | 0.60 | moderate |
| 6 | 0.30 | difficult | 16 | 0.63 | moderate | 26 | 0.68 | moderate |
| 7 | 0.66 | moderate | 17 | 0.58 | moderate | 27 | 0.66 | moderate |
| 8 | 0.72 | easy | 18 | 0.69 | moderate | 28 | 0.55 | moderate |
| 9 | 0.64 | moderate | 19 | 0.50 | moderate | 29 | 0.68 | moderate |
| 10 | 0.54 | moderate | 20 | 0.55 | moderate | 30 | 0.70 | easy |

## Distractor effectiveness (D)

The results of the distractor effectiveness calculation show that 16 questions have ineffective distractors, as the students who chose those distractors did not constitute 5% of the total test takers. Based on the analysis of Tables 1, 2, and 3, it is concluded that questions 5 and 24 were not selected for development into a four-tier diagnostic instrument. The four-tier diagnostic instrument was developed based on 28 open-ended multiple-choice questions.

**Table 3.** The percentage of the distractor effectiveness of short answer questions

| No<br>Opt | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|-----------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| A | 14.6 | 69.4 | 43.0 | 15.4 | 0.0 | 26.7 | 3,2 | 71.8 | 12.6 | 17.8 | 10.8 | 27.8 | 6.4 | 27.3 | 25.8 |
| B | 12.5 | 13.3 | 17.2 | 50.5 | 2.0 | 16.7 | 67,7 | 15.6 | 7.4 | 18.9 | 34.9 | 18.9 | 89.4 | 60.2 | 16.1 |
| C | 14.6 | 8.20 | 26.9 | 26.4 | 94.8 | 24.4 | 3,2 | 1.0 | 63,5 | 5.6 | 12.0 | 4.4 | 4.2 | 9.0 | 53.8 |
| D | 58.3 | 6.20 | 13.0 | 13.2 | 3.1 | 32.2 | 25,8 | 11.5 | 15.7 | 57.8 | 42.2 | 48.9 | 0.0 | 3.4 | 7.5 |

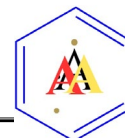| No<br>Opt | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 |
|-----------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| A | 5.7 | 10.2 | 15.9 | 52.7 | 59.5 | 16.3 | 75.6 | 1.0 | 22.2 | 11.5 | 73.9 | 4.5 | 12.8 | 4.3 | 12.8 |
| B | 24.1 | 23.9 | 75.0 | 18.7 | 32.6 | 23.9 | 18.9 | 66.3 | 25.9 | 0.0 | 10.2 | 4.5 | 19.1 | 12.9 | 77.9 |
| C | 68.9 | 7.9 | 7.9 | 14.3 | 6.7 | 54.3 | 3.3 | 30.4 | 45.7 | 60.4 | 11.4 | 20.2 | 56.4 | 12.9 | 2.3 |
| D | 1.1 | 63.6 | 1.1 | 14.3 | 0.0 | 5.4 | 2.2 | 2.1 | 6.1 | 23.0 | 4.5 | 70.8 | 11.7 | 69.9 | 6.9 |

## Four-tier instrument

### Content Validity

The average percentage of instrument feasibility obtained based on content validation was 89.45%. According to Arikunto's (2015: 89) criteria for feasibility levels, the four-tier diagnostic instrument developed by the researcher met the very feasible criteria, so no significant revisions were needed. The four-tier diagnostic instrument was only partially revised in response to suggestions from the validators prior to testing.

### Empirical validity

*Reliability.* Reliability for the B tier (0.858) is higher than for the A tier (0.864) and R tier (0.775). Based on the analysis, the reliability level for the A tier is very high, and for the R tier, it is high. Meanwhile, the reliability level for the entire test (B tier) is very high.

*Difficulty Level.* Based on the average, the developed instrument is moderately difficult. The average of the difficulty index for the A tier (0.53) is lower than that for the R tier (0.55), indicating that more students chose the correct option for the reason (R tier)

than for the answer (A tier). This suggests that most students understand the concept well. Meanwhile, the difficulty index for the B tier (0.42) is lower than that for the A and R tiers because, to answer correctly, students must have a good understanding.

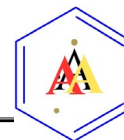**Table 4.** Difficulty Level of A, R, and B tiers.

| No. | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A tier | 0.65 | 0.80 | 0.59 | 0.42 | 0.56 | 0.65 | 0.66 | 0.51 | 0.65 | 0.38 | 0.27 | 0.56 | 0.69 | 0.73 |
| R tier | 0.66 | 0.72 | 0.61 | 0.44 | 0.75 | 0.61 | 0.54 | 0.72 | 0.49 | 0.34 | 0.30 | 0.41 | 0.58 | 0.76 |
| B tier | 0.59 | 0.70 | 0.51 | 0.41 | 0.49 | 0.55 | 0.45 | 0.46 | 0.34 | 0.28 | 0.15 | 0.34 | 0.52 | 0.69 |
| **No** | **15** | **16** | **17** | **18** | **19** | **20** | **21** | **22** | **23** | **24** | **25** | **26** | **27** | **28** |
| A tier | 0.27 | 0.21 | 0.23 | 0.66 | 0.46 | 0.46 | 0.52 | 0.38 | 0.61 | 0.56 | 0.61 | 0.55 | 0.63 | 0.63 |
| R tier | 0.45 | 0.41 | 0.39 | 0.55 | 0.51 | 0.75 | 0.44 | 0.31 | 0.49 | 0.35 | 0.66 | 0.55 | 0.83 | 0.73 |
| B tier | 0.15 | 0.17 | 0.17 | 0.49 | 0.41 | 0.42 | 0.35 | 0.25 | 0.41 | 0.34 | 0.51 | 0.45 | 0.62 | 0.58 |

*Discriminatory Indices.* The analysis results show that the DI for the A, R, and B tiers ranged from poor to good, with no test items having a very good DI. Items 12 and 16 each had a negative DI value of -0.10 and -0.04, respectively. This indicates that the questions cannot distinguish between students with good conceptual understanding and those with low conceptual understanding, so the questions need to be revised. However, there are several considerations before making revisions. In some cases, items with low DI values can be retained because the primary purpose for developing the items was to identify students' conceptual understanding, not to differentiate between high-achieving and low-achieving students (Habiddin & Page, 2019).

**Table 5.** Discriminatory indices of A, R, and B tiers using Pearson Correlation

| No | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A tier | 0,49 | 0,40 | 0,38 | 0,44 | 0,43 | 0,38 | 0,57 | 0,49 | 0,21 | 0,36 | 0,25 | 0,38 | 0,52 | 0,43 |
| R tier | 0,52 | 0,46 | 0,41 | 0,30 | 0,51 | 0,18 | 0,27 | 0,29 | 0,30 | 0,27 | 0,08 | -0,10 | 0,69 | 0,43 |
| B tier | 0,49 | 0,54 | 0,49 | 0,30 | 0,47 | 0,24 | 0,38 | 0,47 | 0,27 | 0,27 | 0,25 | 0,10 | 0,69 | 0,52 |
| **No** | **15** | **16** | **17** | **18** | **19** | **20** | **21** | **22** | **23** | **24** | **25** | **26** | **27** | **28** |
| A tier | 0,19 | 0,13 | 0,16 | 0,57 | 0,13 | 0,35 | 0,24 | 0,02 | 0,41 | 0,43 | 0,29 | 0,41 | 0,40 | 0,40 |
| R tier | 0,21 | -0,04 | 0,33 | 0,52 | 0,15 | 0,35 | 0,35 | 0,10 | 0,47 | 0,07 | 0,35 | 0,41 | 0,12 | 0,20 |
| B tier | 0,25 | 0,05 | 0,16 | 0,58 | 0,13 | 0,33 | 0,36 | 0,16 | 0,52 | 0,16 | 0,44 | 0,49 | 0,43 | 0,29 |

*Distractor Effectiveness.* Based on the analysis results, most distractors are effective, as 84.5% were chosen by more than 5% of test participants.

**Table 6.** Distractor Effectiveness for each option

| No | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | | 7 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Opt | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier |
| A | 8.45 | 66.20 | 80.28 | 2.82 | 14.08 | 22.54 | 5.63 | 5.63 | 9.86 | 7.04 | 26.76 | 18.31 | 12.68 | 11.27 |
| B | 19.72 | 7.04 | 8.45 | 71.83 | 59.15 | 14.08 | 42.25 | 8.45 | 14.08 | 4.23 | 4.23 | 1831 | 11.27 | 19.72 |
| C | 7.04 | 15.49 | 5.63 | 11.27 | 12.68 | 60.56 | 4.23 | 43.66 | 19.72 | 74.65 | 4.23 | 2.82 | 66.20 | 53.52 |
| D | 64.79 | 11.27 | 5.63 | 14.08 | 14.08 | 2.82 | 47.89 | 42.25 | 56.34 | 14.08 | 64.79 | 60.56 | 9.86 | 15.49 |

| No | 8 | | 9 | | 10 | | 11 | | 12 | | 13 | | 14 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Opt | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier |
| A | 50.70 | 11.27 | 19.72 | 30.99 | 38.03 | 7.04 | 21.13 | 29.58 | 25.35 | 33.80 | 14.08 | 14.08 | 73.24 | 11.27 |
| B | 15.49 | 4.23 | 7.04 | 49.30 | 42.25 | 33.80 | 33.80 | 14.08 | 56.34 | 40.85 | 12.68 | 14.08 | 16.90 | 76.06 |
| C | 19.72 | 12.68 | 64.79 | 12.68 | 8.45 | 16.90 | 18.31 | 14.08 | 12.68 | 18.31 | 69.01 | 14.08 | 4.23 | 1.41 |
| D | 14.08 | 71.83 | 8.45 | 7.04 | 11.27 | 42.25 | 26.76 | 42.25 | 5.63 | 7.04 | 4.23 | 57.75 | 5.63 | 11.27 |

| No | 15 | | 16 | | 17 | | 18 | | 19 | | 20 | | 21 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Opt | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier |
| A | 26.76 | 45.07 | 32.39 | 25.35 | 18.31 | 14.08 | 19.72 | 53.52 | 11.27 | 49.30 | 7.04 | 74.65 | 15.49 | 16.90 |
| B | 2254 | 12.68 | 25.35 | 23.94 | 40.85 | 39.44 | 66.20 | 14.08 | 25.35 | 26.76 | 16.90 | 9.86 | 12.68 | 16.90 |
| C | 28.17 | 18.31 | 21.13 | 40.85 | 22.54 | 26.76 | 5.63 | 16.90 | 46.48 | 15.49 | 29.58 | 12.68 | 52.11 | 43.66 |
| D | 22.54 | 23.94 | 21.13 | 9.86 | 18.31 | 19.72 | 8.45 | 14.08 | 45.07 | 7.04 | 46.48 | 2.82 | 19.72 | 22.54 |

| Soal | 22 | | 23 | | 24 | | 25 | | 26 | | 27 | | 28 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Opsi | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier | A tier | R tier |
| A | 53.52 | 49.30 | 60.56 | 21.13 | 4.23 | 35.21 | 60.56 | 12.68 | 14.08 | 54.93 | 2.82 | 4.23 | 63.38 | 9.86 |
| B | 38.03 | 30.99 | 11.27 | 18.31 | 23.94 | 39.44 | 19.72 | 66.20 | 16.90 | 16.90 | 28.17 | 83.10 | 9.86 | 11.27 |
| C | 5.63 | 14.08 | 18.31 | 11.27 | 56.34 | 15.49 | 9.86 | 14.08 | 14.08 | 12.68 | 63.38 | 8.45 | 22.54 | 73.24 |
| D | 2.82 | 5.63 | 9.86 | 49.30 | 15.49 | 9.86 | 9.86 | 7.04 | 54.93 | 15.49 | 5.63 | 4.23 | 4.23 | 5.63 |

*Validity.* The analysis results show that most of the developed questions are valid, but some items are not. A total of 3 questions were invalid at the A tier, 5 questions were invalid at the R tier, and 4 questions were invalid at the B tier. These invalid questions need to be considered for revision based on other parameters, namely difficulty level, discrimination index, and distractor effectiveness. Based on empirical validation analysis, item 16 was discarded, items 9, 10, 11, 12, and 19 were retained with revisions, and items 15 and 22 did not require revision.

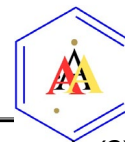**Table 7.** Validity of A, R, and B tiers

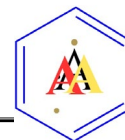| No | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A tier | $r_{xy}$ | 0.483 | 0,.01 | 0.399 | 0.456 | 0.497 | 0.488 | 0.621 | 0.544 | 0.367 | 0.388 | 0.417 | 0.429 | 0.615 | 0.628 |
| R tier | $r_{xy}$ | 0.624 | 0.611 | 0.577 | 0.329 | 0.608 | 0.308 | 0.377 | 0.373 | 0.229 | 0.389 | 0.119 | -0.140 | 0.660 | 0.452 |
| B tier | $r_{xy}$ | 0.582 | 0.603 | 0.628 | 0.417 | 0.613 | 0.407 | 0.547 | 0.556 | 0.181 | 0.421 | 0.344 | 0.121 | 0.728 | 0.645 |
| No | | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 |
| A tier | $r_{xy}$ | 0.286 | 0.179 | 0.400 | 0.641 | 0.216 | 0.424 | 0.295 | 0.106 | 0.511 | 0.560 | 0.447 | 0.476 | 0.484 | 0.469 |
| R tier | $r_{xy}$ | 0.224 | 0.044 | 0.484 | 0.574 | 0.233 | 0.329 | 0.483 | 0.231 | 0.517 | 0.283 | 0.481 | 0.512 | 0.232 | 0.262 |
| B tier | $r_{xy}$ | 0.442 | 0.196 | 0.371 | 0.646 | 0.274 | 0.331 | 0.534 | 0.201 | 0.561 | 0.320 | 0.576 | 0.509 | 0.475 | 0.407 |

## CONCLUSIONS

The resulting product development is a four-tier diagnostic instrument to identify misconceptions of 11th-grade science students regarding salt hydrolysis material. The developed instrument consists of 27 questions. The specifications of the resulting product are: (1) The developed four-tier diagnostic instrument consists of four tiers, with the first tier being questions and answers with four answer options, the second tier representing students' confidence level in choosing the first tier on a scale of 1-5 (1=guessing only; 2=not sure; 3=moderate; 4=sure; 5=very sure), the third tier being the reason for the first tier, and the third tier being the confidence level for the reason

on a scale of 1-5 (1=guessing only; 2=not sure; 3=moderate; 4=sure; 5=very sure); (2) The reason choices used are based on students' reasons in the open-ended multiple-choice initial test and relevant literature; (3) The developed four-tier diagnostic instrument consists of at least one question per indicator; (4) The instructions for answering the questions in the developed instrument include general instructions for answering the presented questions.

## REFERENCES

Amala, F., & Habiddin, H. (2022). Pemahaman konsep dalam topik sifat asam basa larutan garam: studi pada siswa SMA di Blitar. *Jurnal Zarah*, *10*(2), 91–100. https://doi.org/10.31629/Zarah.V10I2.4321

Amelia, D., Marheni, M., & Nurbaity, N. (2014). Analisis Miskonsepsi Siswa Pada Materi Hidrolisis Garam Menggunakan Teknik CRI (Certainty Of Response Index) Termodifikasi. *JRPK: Jurnal Riset Pendidikan Kimia*. https://doi.org/10.21009/jrpk.041.05

Ardina, D., & Habiddin, H. (2023). Acid-base properties of salt solution: Study at a secondary school in Banyuwangi. *AIP Conference Proceedings*, *2569*(1), 30018. https://doi.org/10.1063/5.0112074

Beichner, R. J. (1994). Testing student interpretation of kinematics graphs. *American Journal of Physics*, *62*(8), 750–762. https://doi.org/10.1119/1.17449

Billah, A. A., Octavia, D. A., Sholihah, F. H., Febriliyanto, M. R., Aisyah, N. N., Rohmah, N. H., Rahayu, S. M. S., Salsabela, T. A., Basimin, M. Q., & Habiddin, H. (2024). University Students' Understanding of Intermolecular Forces: A Comparison of First & Third-Year Students. *STEM Education International*, *1*(1 SE-Research Articles), 1–7. https://sciencesustain.com/index.php/stem/article/view/5

Demircioğlu, G., Ayas, A., & Demircioğlu, H. (2005). Conceptual change achieved through a new teaching program on acids and bases. *Chemistry Education Research and Practice*, *6*(1), 36–51. https://doi.org/10.1039/B4RP90003K

Gurel, D. K., Eryilmaz, A., & McDermott, L. C. (2015). A review and comparison of diagnostic instruments to identify students' misconceptions in science. *Eurasia Journal of Mathematics, Science and Technology Education*, *11*(5), 989–1008. https://doi.org/10.12973/eurasia.2015.1369a

Gurel, D. K., Eryilmaz, A., & McDermott, L. C. (2017). Development and application of a four-tier test to assess pre-service physics teachers' misconceptions about geometrical optics. *Research in Science & Technological Education*, *35*(2), 238–260. https://doi.org/10.1080/02635143.2017.1310094

Habiddin, H., & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, *19*(3), 720–736. https://doi.org/10.22146/ijc.39218

Habiddin, H., & Page, E. M. (2023). Uncovering Students' Genuine Misconceptions: Evidence to Inform the Teaching of Chemical Kinetics. *Acta Chimica Slovenica*, 184–195.

Laliyo, L. A. R., Hamdi, S., Pikoli, M., Abdullah, R., & Panigoro, C. (2021). Implementation of four-tier multiple-choice instruments based on the partial credit model in evaluating students' learning progress. *European Journal of Educational Research*, *10*(2), 825–840. https://doi.org/10.12973/EU-JER.10.2.825

Maratusholihah, N. F., Sri, R., & Fauziatul, F. (2017). Analisis Miskonsepsi Siswa SMA pada Materi Hidrolisis Garam dan Larutan Penyangga. *Jurnal Pendidikan:Teori, Penelitian, Dan Pengembangan*, *2*(7), 919—926.

Novak, J. D. (1990). Concept Mapping - A Useful Tool For Science-Education. *Journal of*

*Research in Science Teaching*, *27*(10), 937–949.

Orwat, K., Bernard, P., & Mikuli, A. M. (2017). Alternative Conceptions of Common Salt Hydrolysis Among Upper-Secondary-School Students. *Journal of Baltic Science Education*, *16*(1), 64–76. https://doi.org/https://doi.org/10.33225/jbse/17.16.64

Osborne, R. ., & Gilbert, J. K. (1980). A Method for Investigating Concept Understanding in Science. *European Journal of Science Education*, *2*(3), 311–321.

Pesman, H., & Eryilmaz, A. (2010). Development of a Three-Tier Test to Assess Misconceptions About Simple Electric Circuits. *Journal of Educational Research*, *103*(3), 208–222.

Seçken, N. (2010). Identifying student's misconceptions about salt. *Procedia - Social and Behavioral Sciences*. https://doi.org/10.1016/j.sbspro.2010.03.004

Taber, K. S. (1999). Ideas About Ionisation Energy: A Diagnostic Instrument. *School Science Review*, *81*(295), 97–104.

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.

Tuysuz, C. (2009). Development of two-tier diagnostic instrument and assess students' understanding in chemistry. *Scientific Research and Essays*, *4*(6), 626–631.

Ulfah, M., Ulfah, M., Erlina, E., Pratiwiningrum, F. M., Wafiq, A. F., & Juahir, Y. B. (2024). Misconceptions of Intermolecular Forces in General Chemistry Courses. *Jurnal IPA & Pembelajaran IPA*, *8*(1), 39–51. https://doi.org/10.24815/jipi.v8i1.35641