# Integration of Google Sheets with Deep Learning-based Evaluation Methods to Measure the Analytical Skills of PPG Students in Creating HOTS Questions

# Integrasi *Google Sheet* dengan Metode Evaluasi berbasis *Deep Learning* untuk Mengukur Kemampuan Analitis Mahasiswa PPG dalam Pembuatan Soal HOTS

Adip Wahyudi[1*], Slamet Arifin[1], Evi Fitriana[1], Feri Fahrian Maulana[2], Dynda Prista[1],
Muhammad Aris Firmansyah[3], Radya Wirawan Nurjati[4]
[1] Universitas Negeri Malang, Malang, Indonesia
[2] Universitas Brawijaya, Malang, Indonesia
[3] PT. Alfahuma Rekayasa Teknologi, Malang, Indonesia
[4] National Yang Ming Chiao Tung University, Hsinchu City, Taiwan
[*] E-mail: adip.wahyudi.fis@um.ac.id

*Abstract*

*21st-century education requires the use of creative interactive tools to make learning outcome evaluation more effective. The main challenge is low student participation due to conventional assessment practices. The Google Sheet evaluation tool that integrates deep learning offers a solution to the problem of learning evaluation. This study focuses on integrating an evaluation tool into Google Sheets using deep learning in the Learning Evaluation course in the PPG program and on testing its impact on learning outcomes. The research method used R&D with the ADDIE model, with data collected through a pretest-posttest design across two groups: the experimental group that used deep learning-based Google Sheets evaluation tools and the control group that used a conventional approach. After normality and homogeneity tests confirmed that the statistical prerequisites were met, the analysis continued with a paired t-test. The results showed significance (p < 0.05), indicating that the deep learning-based Google Sheets evaluation tool was effective in improving learning outcomes, while also encouraging student motivation and engagement. In conclusion, the deep learning-based Google Sheets evaluation tool is suitable for implementation in the PPG program, with opportunities for further development through the enrichment of collaboration features and adaptive evaluation to strengthen the appeal and impact of learning.*

*Keywords: Google sheets, deep learning, evaluation tools, teacher professional education.*

**Abstrak**

Pendidikan abad 21 menuntut pemanfaatan alat interaktif yang kreatif agar evaluasi hasil belajar lebih efektif. Tantangan utamanya adalah rendahnya partisipasi mahasiswa akibat praktik penilaian yang masih konvensional. Alat evaluasi *google sheet* yang integrasikan *deep learning* menawarkan solusi untuk permasalahan evaluasi pembelajaran. Penelitian berfokus mengintegrasikan alat evaluasi dengan *google sheet* berbasis *deep learning* pada mata kuliah evaluasi pembelajaran di program PPG serta menguji dampaknya terhadap peningkatan hasil belajar. Metode penelitian menggunakan R&D dengan model ADDIE dengan data dikumpulkan lewat desain *pretest posttest* pada dua kelompok, yakni kelompok eksperimen yang menggunakan alat evaluasi *google sheets* berbasis *deep learning* dan kelompok kontrol yang memakai pendekatan konvensional. Setelah uji normalitas dan homogenitas memastikan prasyarat statistik terpenuhi, analisis dilanjutkan menggunakan uji t berpasangan. Hasilnya menunjukkan signifikansi (p < 0,05), menandakan alat evaluasi *google sheet* berbasis *deep learning* efektif meningkatkan capaian belajar, sekaligus mendorong motivasi dan keterlibatan mahasiswa. Kesimpulannya alat evaluasi *google sheets* berbasis *deep learning* layak diimplementasikan pada program PPG, dengan peluang pengembangan berikutnya melalui pengayaan fitur kolaborasi dan evaluasi adaptif untuk memperkuat daya tarik dan dampak pembelajaran.

**Kata Kunci:** *Google sheet, deep learning*, alat evaluasi, pendidikan profesi guru.

## 1. Introduction

The global education paradigm has undergone a profound transformation. Many professions in the future will require advanced analytical skills and creativity, prompting countries such as Finland and Singapore to integrate Higher Order Thinking Skills into their education curricula (Dos Santos et al., 2025). This change is not merely a trend but a strategic necessity for developing competitive human resources (Kacmaz & Dubé, 2025; Maulana et al., 2023). In Indonesia, implementing HOTS faces various challenges. Research conducted by Khoiriyah & Shaleh (2025) shows that only a few schools can implement the Higher Order Thinking Skills curriculum effectively, while the rest still rely on conventional learning methods, especially in 3T (Underdeveloped, Frontier, and Outermost) areas. The imbalance between progressive policies and the actual state of educational infrastructure has led to systemic inefficiencies (Mir et al., 2024).

Education in the 21st century requires the development of Higher Order Thinking Skills (HOTS) as the foundation of learning. Although the Indonesian curriculum has adopted a 21st century learning approach, practice in the field is still hampered by the dominance of evaluation instruments that only measure Lower Order Thinking Skills (LOTS) (Nasbi, 2017). Teachers, as the main agents of education, play a vital role, so their competence needs to be honed through Teacher Professional Education programs (Saputra & Nofrion, 2022). The main focus of this competency development lies in analytical skills, which are conceptually defined as C4 (Analyzing) cognitive abilities according to Bloom's revised taxonomy, including the skills of distinguishing, organizing, and attributing elements of information to construct complex test stimuli. However, there is currently a significant gap between these competency requirements and the available evaluation infrastructure. PPG students are required to reconstruct material into high-level questions, but the evaluation instruments used remain conventional and manual, failing to objectively assess the semantic structure and cognitive level of the questions. As a result, the validity of HOTS question construction is not well verified.

Based on observations in the Geography PPG Program at Universitas Negeri Malang (UM), empirical findings indicate that 65% of 40 students in the program have difficulty formulating analysis and evaluation questions in accordance with Bloom's revised taxonomy. Classroom observations show that 70% of the 40 students' questions remain at the C1-C3 level, indicating a lack of contextualized question structure. Interviews with PPG lecturers reveal that the main problem is the absence of instruments capable of assessing question quality while providing objective feedback. This leads to differences in the quality of the questions produced. Google Sheets has been used to compile questions due to its ease of organizing items and basic statistical analysis (Surahman et al., 2024). However, the previous instrument, which used Google Sheets, lacked features to automatically evaluate the cognitive complexity of questions. In fact, deep learning technology has already been proven in research (Yudistira & Kom, 2024) to be capable of classifying text by cognitive difficulty level using a Natural Language Processing model. The integration of these two tools has the potential to provide a holistic solution, but no research has implemented it in a structured manner in the PPG context.

Deep learning can process structured and unstructured data, such as text, images, and audio, through multi-level representation learning, where each layer of the network transforms the input into increasingly abstract and meaningful representations (Nasser et al., 2025). This capability enables extensive data collection and algorithms that assign different weights to each element of text. Thus, when used in education, Deep learning has demonstrated its effectiveness in tasks such as classifying the cognitive difficulty level of questions, generating higher-order thinking skills questions, and performing automatic feedback analysis (Sumy et al., 2023; Vanhees et al., 2025).

As evidenced by research (Kasanah & Pratama, 2024), which achieved 92% accuracy in categorizing questions according to Bloom's revised taxonomy.

Google Sheets functions as more than just a spreadsheet application; it serves as an intelligent evaluation interface through the use of Google Apps Script customization features (Parra et al., 2021). This feature is a strategic key that enables the integration of Deep Learning into student worksheets. In the context of the Geography Teacher Education Program (PPG) at UM, this integration transforms conventional spreadsheets into an interactive analytical ecosystem. Students not only enter the question text, but the system also automatically processes it to diagnose the cognitive level in real time. This instant feedback mechanism directly supports students' analytical skills in formulating complex questions, making Google Sheets an effective pedagogical tool for validating evaluation instruments (Rebman et al., 2023).

Research has made substantial progress in using Deep Learning for educational assessment, as exemplified by a transformer model that achieved 85% accuracy in categorizing Higher Order Thinking Skills problems (Chen & Han, 2025). However, these programs operate independently and are not integrated with established educator platforms, such as Excel, so they cannot be used for teacher training. Research conducted by Wahyudi et al. (2025) in the PPG Program at Universitas Negeri Malang revealed that 80% of students prefer to use Excel to compile questions. However, this application does not have an automatic cognitive analysis function, and validators must check it offline. Research by Aniek et al. (2022) found that the limitations of traditional tools in providing objective assessments of the quality of Higher Order Thinking Skills questions create a gap between the demands for developing students' analytical skills and current evaluation methodologies. This problem is exacerbated by the prevalence of C1-C3 level questions (LOTS) created by students, as revealed in field observations.

Based on previous research, the researchers filled the research gap by integrating Google Sheets with a bidirectional LSTM that analyzes the linguistic structure of questions, classifies cognitive levels according to Bloom's taxonomy, and offers recommendations to improve higher-order thinking skills. This integration seeks to align technological capabilities with educational needs while improving the objectivity and consistency of the quality of questions prospective educators generate for students in geography lessons.

## 2. Literatur Review

## 2.1. HOTS and the Learning Evaluation Paradigm

This study is based on the Constructivist theory, which views question construction not merely as an administrative task but as a cognitive process in which students construct meaning through authentic experiences (Sugrah, 2020). The main variable measured in this study, namely analytical ability, is based on Anderson & Krathwohl's Taxonomy. The main variable in Bloom's taxonomy has become the standard for content validity, emphasizing the higher-order cognitive processes of analyze, evaluate, and create (Kasanah & Pratama, 2024).

To support the development of these abilities, this study adopts the Assessment for Learning principle. Evaluation does not only serve as a final measurement tool but also as a continuous feedback mechanism that drives improvements in question quality (Misnawati et al., 2025). The integration of this theory confirms that improvements in PPG students' abilities can occur only if the system provides interventions in the form of specific suggestions for improvement on the draft questions they have compiled.

### 2.2. Deep Learning in Automated Assesment

The application of technology in this study is based on the latest literature on Deep Learning and Natural Language Processing. Research (Akdemir & Barışçı, 2024) shows the capacity of Deep Learning models to capture semantic patterns and natural language reasoning structures, which are relevant for detecting whether an item meets HOTS criteria or only LOTS. Specifically, the Deep Learning architecture in this system serves as an intelligent validator, replacing the slow, manual validation process. According to Gardner et al. (2021), automated systems offer advantages in assessment consistency. The inference pipeline, from text preprocessing to result visualization, produces actionable diagnostic analytics (Dagal et al., 2025). This converts raw scores into qualitative feedback that guides students in revising their stimuli and answer options.

In conclusion, Deep Learning in this study serves as an intelligent validator that automatically distinguishes the quality of HOTS and LOTS questions through language pattern analysis. Its existence overcomes the weaknesses of manual validation by offering greater speed and consistency in assessment. More than just giving scores, this system provides concrete feedback for improvement, directly guiding students to refine the quality of the questions they create.

### 2.3. Technology Integration and Cognitive Load

The selection of Google Sheets as the main interface is justified by Cognitive Load Theory (CLT). Research (Gkintoni et al., 2025) emphasizes the importance of minimizing extraneous cognitive load so that users' mental energy is focused on germane load. Using a familiar spreadsheet platform minimizes technical barriers for students, allowing them to focus on the substance of question development. The acceptance of this technology was also analyzed through the "Technological Pedagogical Content Knowledge" and "Technology Acceptance Model" frameworks. Technological Pedagogical Content Knowledge explains the synergy among geography content, pedagogical evaluation, and AI technology, while the Technology Acceptance Model predicts that students' perceptions of Google Sheets' ease of use will increase their adoption of the system (Huang et al., 2025). This interactive format has been empirically proven to increase engagement and reduce boredom in complex learning processes.

In conclusion, using Google Sheets aims to free students from the mental energy required to learn complex new applications, allowing them to focus fully on the substance of question preparation. Because this platform is already very familiar in their daily lives, the advanced technology embedded in it feels light and easy to accept, with few significant technical barriers. This familiar approach ultimately makes the learning process more lively, less boring, and able to increase student participation.

### 2.4. Conceptual Framework

Based on the synthesis of the above theories, this study's conceptual framework links problems, technological solutions, and expected outcomes. The process begins with student input by creating draft questions that are often still at the LOTS level. Through a system intervention that integrates Deep Learning algorithms into the Google Sheet interface, the system performs real-time semantic structure analysis. The system's output is not only classification labels, but also constructive feedback that triggers students to make independent revisions. This iterative cycle gradually builds students' understanding, ultimately leading to an increase in the study's dependent variable: analytical skills in developing HOTS-based evaluation instruments.

## 3. Research Method

The Google Sheets integration system with a deep learning-based evaluation method was developed using a user-centered ADDIE approach in this product. The researchers chose this approach because it can overcome complex challenges by putting themselves in the position of users, namely PPG students and teachers, iterating quickly on prototypes, and producing creative solutions that can be applied in geography classes. The five adaptive stages of the ADDIE process, as described by Maxnun et al. (2024), are analysis, design, development, implementation, and evaluation, as illustrated in Figure 1.
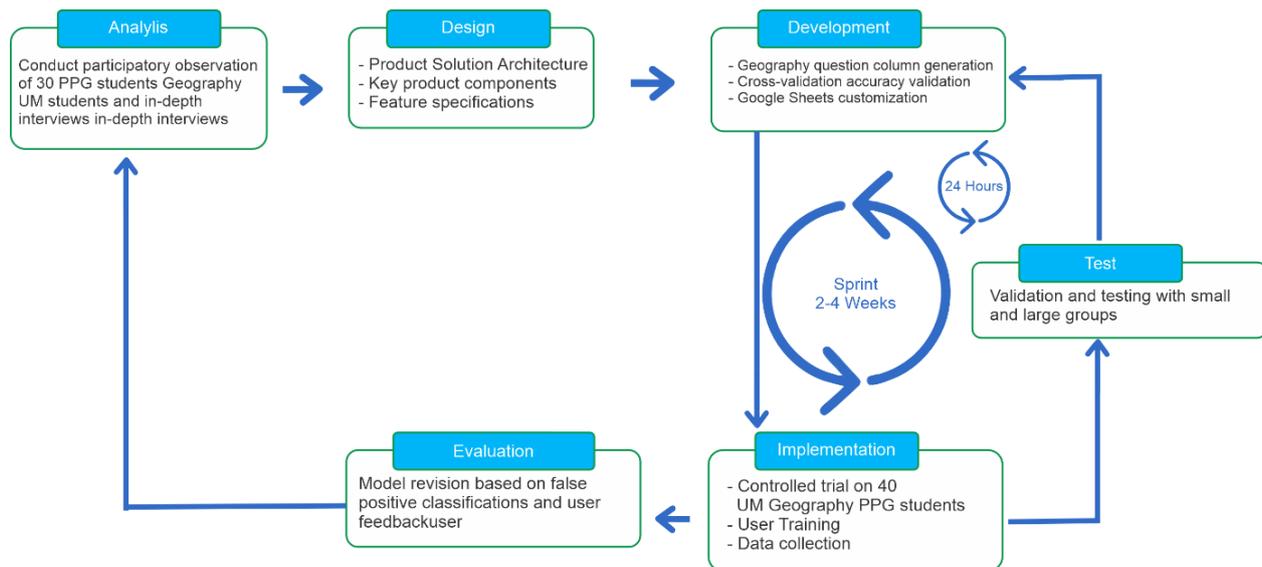


**Figure.1**. ADDIE Research Flow

The subjects of this study involved 40 students of the Teacher Professional Education Program (PPG) at Universitas Negeri Malang in the 2025 academic year. Sampling was conducted using Purposive Sampling, with respondents characterized as PPG Geography Education students enrolled in the Learning Evaluation course. To test the system's effectiveness, respondents were divided into two groups in a quasi-experimental design: 20 students in the experimental group using Google Sheets integrated with Deep Learning, and 20 students in the control group using the conventional Excel method.

Qualitative and quantitative data were obtained directly from graduate students in the PPG program at Universitas Negeri Malang through three structured instruments (Ardiansyah et al., 2023; Nasution, 2023). First, an expert validation questionnaire to evaluate content validity and practicality. Second, a user response questionnaire trial for education practitioners and students using a Likert scale to assess Content Quality, Interactive Design and Affordability, Implementation in Learning, and Practical Needs. Third, implementation with students to measure the improvement in their skills in formulating questions based on Higher Order Thinking Skills (HOTS) before the pre-test and after the post-test. The open instrument produced qualitative data in the form of criticism and suggestions, while the closed instrument produced validity scores and learning outcomes, which were then analyzed for content validity using Aiken's V Formula:

$$V = \frac{\sum \left( \frac{r_i - r_{min}}{r_{max} - r_{min}} \right)}{n}$$

Explanation:

$V$ : Validity coefficient

$r_i$ : Score given by validator i

$r_{min}$ : Lowest score on the rating scale

$r_{max}$ : Highest score on the rating scale

$n$ : Number of validators

Validity criteria $V \geq 0.75$ (Ligasan & Doysabas, 2024), (b) quantitative statistical evaluation including Cronbach's Alpha reliability test ($\alpha \geq 0.7$) and paired sample t-test ($p < 0.05$), and (c) thematic analysis (Bujang et al., 2018) for qualitative data, facilitated by SPSS 26.0 software. The scores obtained from respondents' evaluations are quantitative data that will be converted to qualitative data by referring to the assessment range in Table 1 as follows.

**Table.1**. Content Criteria and Suitability

| No | Description | Class |
|----|-------------|-------|
| 1 | 0%-20% | Very Unsuitable (STS) |
| 2 | 21%-40% | Less Suitable (KL) |
| 3 | 41%-60% | Undecided (R) |
| 4 | 61%-80% | Suitable (L) |
| 5 | 81%-100% | Very Good (SL) |

Source: (Sumartini et al., 2020)

This study used the pretest and posttest scores of graduate students in the Geography PPG program at Malang State University for data analysis. Before conducting a paired-samples t-test, normality and homogeneity tests were performed to ensure the data met the necessary statistical assumptions. The pretest and posttest results were then evaluated using a paired sample t-test to test the research hypothesis. As a result, the t-test results were compared with the predetermined significance threshold.

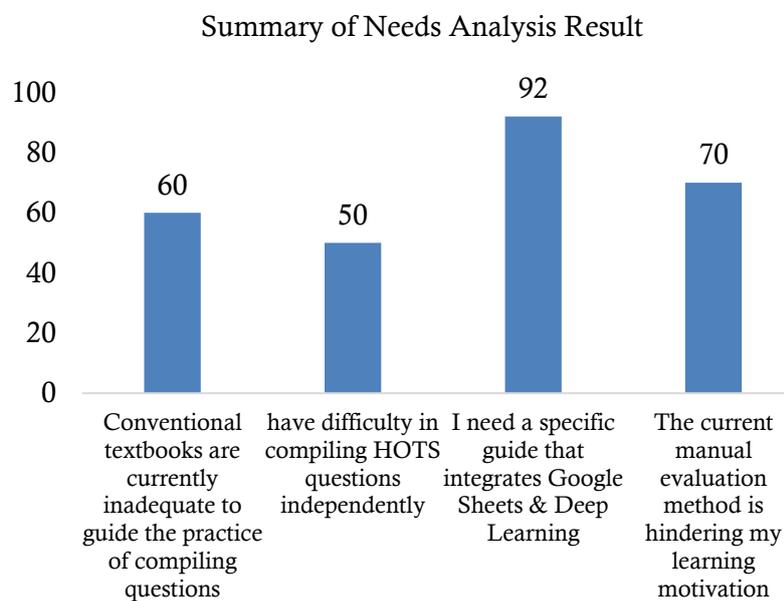## 4. Results and Discussion

### 4.1. Research Results

#### 4.1.1. Analysis

Interviews with geography PPG lecturers revealed that learning evaluation courses still use conventional pedagogical approaches centered on instructors, resulting in decreased engagement and enthusiasm among geography students. Significant findings show that students have difficulty understanding the material in the learning outcome evaluation module and show a lack of interest in static media such as PowerPoint. This study provides a unique teaching tool, namely an evaluation tool that integrates deep learning-based Google Sheets to overcome these difficulties. This solution aims to increase learning enthusiasm and reduce boredom, while also serving as a means to assess the analytical abilities of PPG Geography students in formulating Higher Order

Thinking Skills questions objectively, thereby promoting a comprehensive understanding of evaluation concepts.

A user needs analysis was conducted to identify the evaluation tools needed by graduate students in the geography program (PPG) at Universitas Negeri Malang. A Google Form questionnaire was distributed to geography students at Malang State University to facilitate data collection. Based on the results of the needs analysis in Figure 2, the need to develop evaluation tools was confirmed by the significant gap between competency requirements and the availability of learning facilities. It was noted that 60% of students stated that conventional learning resources were no longer adequate to guide instrument development, and 50% of respondents experienced significant difficulties in constructing HOTS questions independently. Students felt that the old methods reduced their learning motivation by 70%. As a solution, the majority of respondents (92%) specifically wanted the integration of Google Sheets technology based on Deep Learning.



**Figure.2**. Results of Needs Analysis

### 4.1.2. Design

A Deep Learning-based Google Sheets evaluation tool was created to improve the learning experience of PPG students in formulating geography questions oriented towards Higher Order Thinking Skills. This process began with selecting an optimized platform as an effective educational tool. This involved designing a storyboard as a framework for the system interface, from the home page to the evaluation module, as well as developing content through the identification of core competencies, determination of material types, and selection of reference learning sources. The result is an integrated Deep Learning evaluation tool aligned with the research objectives that assesses the analytical skills of PPG students through a combination of Google Sheets and Deep Learning.

### 4.1.3. Development

The development stage of the Deep Learning-based Google Sheets evaluation tool focuses on creating and testing the feasibility of an interactive evaluation system. The process begins with formulating learning objectives and preparing materials in accordance with the syllabus. Next, the evaluation tool is developed based on the initial design, complete with a user manual for students.

The evaluation tool has several features, including review cards, introduction, cover, data entry, results analysis, deep question item analysis, remedial enrichment, remedial attendance list, and STEM (Science, Technology, Engineering, and Mathematics) deep-based remedial material. The final result is a structured evaluation tool based on competencies, complete with a questionnaire sheet to measure the validity of the user response tool.

The feature details are illustrated in Figure 3, which shows the customized review card sheet for the high school geography review card that has been created. Figure 4 presents the difficulty-level analysis, which aims to determine whether the question is easy, medium, or difficult. Figure 5 shows the cover page. Create the cover page according to your creativity. Figure 6 shows the analysis of questions tested in previous exams or other assignments. This is the result of an analysis of multiple-choice questions based on deep learning. It displays the identity of the question analysis. Fill in the form as instructed, and also fill in the answer key in the column marked with a red circle to make it structured. Figure 7 shows the results of the question analysis. Green indicates the top group with high scores, while red indicates the bottom group with low scores. Figure 8 shows the results of the question difficulty analysis. We can see which questions are difficult and which are not, and whether they have been revised or rejected/accepted. Figure 9 shows the remedial and enrichment display, where you can see the names of students undergoing remediation along with the suggestions they must follow according to the predetermined rules. Figure 10 shows the remedial attendance list, where students decide whether to do remedial work together or individually. However, they must still pay attention to the teacher's schedule, especially for geography lessons. Figure 11, the last page, is the remedial form. This form can be adjusted to each student's remedial needs. This form has been integrated with Google Gemini Artificial Intelligence for remedial recommendations based on Deep Learning material.



**Figure.3**. Initial Display of the Review Card



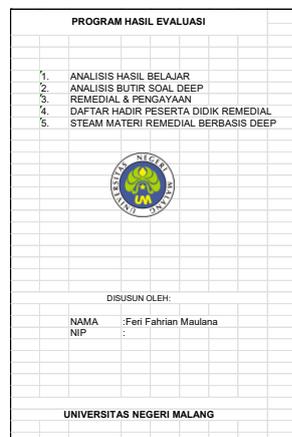**Figure.4**. Introduction to Item Analysis



**Figure.5**. Cover Display



**Figure.6**. Analysis of Questions and Objective Items

**Figure.7**. Item Analysis in One Class



**Figure.8**. Results of Objective Item Analysis



**Figure.9**. Remedial & Enrichment



**Figure.10**. Remedial Attendance List



**Figure.11**. Deep AI Remedial Material

The content validation process involved two subject matter experts who assessed three main indicators, namely curriculum suitability, depth of concept, and reference currency. Based on Figure 12, the post-revision material aspect recorded an average percentage score of 90%. According to the Likert scale feasibility interpretation standard in Table 1, this score is categorized

9

as 'Very Feasible' and does not require major revisions. This high validity indicates that the material has been explicitly, accurately, and relevantly compiled to meet the competency needs of pre-service geography teacher education students. This is crucial, as unstructured material can trigger excessive abstraction and hinder understanding retention (Setiawan et al., 2023). These findings empirically validate that the material is ready for limited-scale testing with practitioners and students, although some minor refinements are still recommended before full adoption.
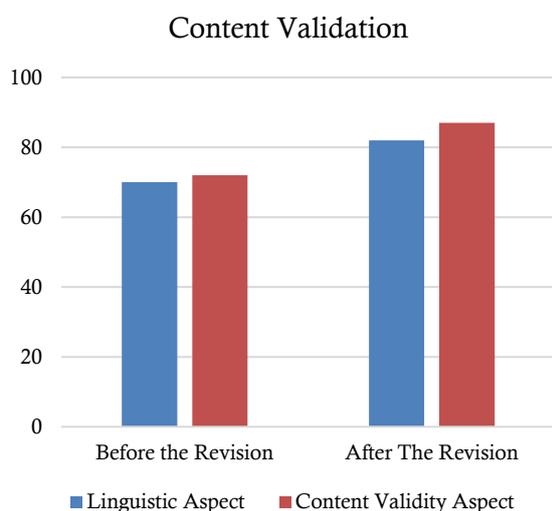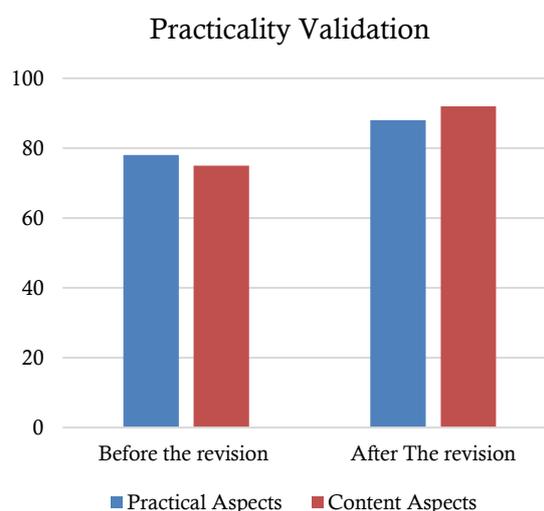


**Figure.12**. Material Validation

**Figure.13**. Practicality Validation

Based on the data visualization in Figure 13, the practicality test of this Deep Learning-based evaluation tool resulted in an average score of 88.5%. When converted according to the Likert scale guidelines in Table 1, this achievement falls into the Very Practical category. This high positive response indicates three main indicators, namely user-friendliness, AI analysis feature navigation flow, and system stability, which have been tested and declared very feasible to proceed to the field trial stage with a wider range of respondents. This technology provides a more immersive learning experience through visualization, helping students understand evaluation information through relevant images. This achievement is in line with the cone of experience theory (Dale, 1969), which emphasizes the importance of multisensory learning in improving teaching and learning methods.
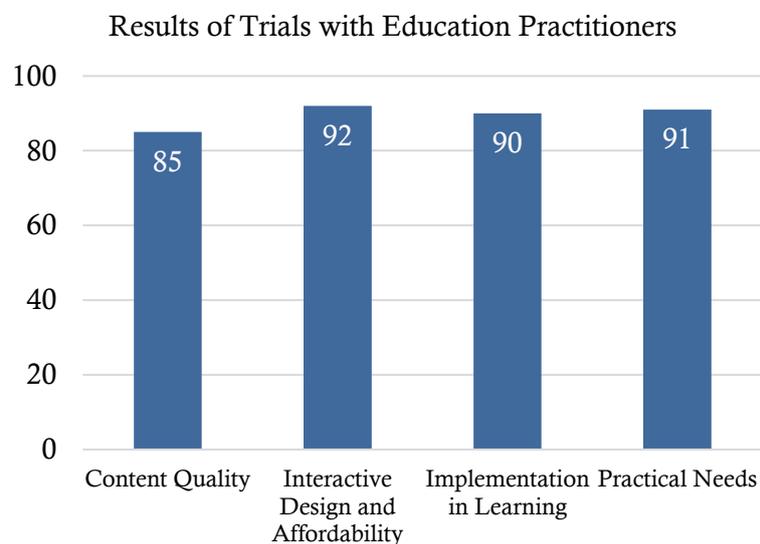
The use of deep learning-based evaluation tools is expected to improve the experience of Geography PPG students in creating Higher Order Thinking Skills questions. Evaluation tools enable students to understand the process of question creation in depth (Sandika et al., 2025). This tool also facilitates students' understanding and development of skills in creating Higher Order Thinking Skills questions, enabling them to create deep, learning-based questions for high school students. This tool has been verified and is ready for testing with education practitioners and students of the Geography Education Master's Program (PPG).

### 4.1.4. Implementation

After undergoing expert validation during the development stage, a deep learning-based evaluation tool integrated with Google Sheets was applied to users to assess the analytical abilities of PPG students in compiling Higher Order Thinking Skills questions. The application involved lecturers as education practitioners and 20 postgraduate PPG students as the experimental group. The trial by practitioners focused on ease of use, the relevance of the model's analytical feedback,

and educators' acceptance of the Google Sheet-based workflow. The assessment was conducted using a questionnaire containing 10 statements accompanied by comment and suggestion columns. The aspects assessed included language clarity, the suitability of system components such as Google Sheet templates, the mechanism for uploading Higher Order Thinking Skills items, the deep learning model inference results dashboard, the regularity of the material structure and assessment flow, and the overall quality of presentation and reporting.
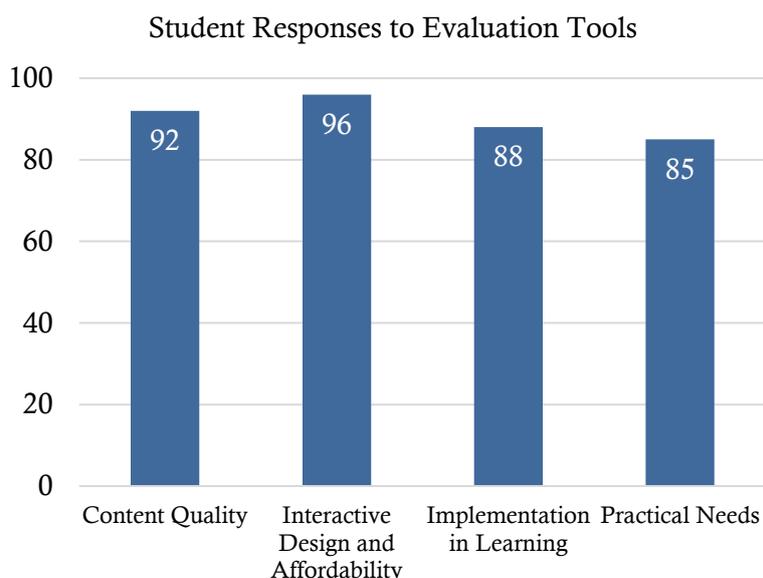
This stage is intended to assess student responses to a deep learning-based Google Sheet-integrated evaluation in PPG learning, specifically whether the system can improve students' understanding of learning evaluation and prepare Higher Order Thinking Skills questions for students. The tested material emphasizes interactive visual presentations that are relevant to students' needs, ranging from question-item input templates in Google Sheets to analytical dashboards of model-prediction results, thereby making the learning process more concise and focused. This implementation not only deepens students' conceptual understanding but also strengthens lecturers' roles as facilitators who guide the interpretation of automatic feedback, while fostering students' learning independence in accordance with constructivist principles.



**Figure.14**. Graph of Test Results for Education Practitioners

Figure 14 shows the results of trials conducted by education practitioners, with an average feasibility of 95%, placing the deep learning-based Google Sheet-integrated evaluation tool in the highly feasible category. These findings confirm that the solution developed, ranging from the HOTS question input template in Google Sheet to the analytical dashboard of the model results, meets the quality criteria to proceed to a more extensive testing phase and is ready to be implemented for PPG students.

Data collection was carried out during offline learning in three sessions guided by PPG lecturers. Before the trial, students were directed to access the deep learning-based Google Sheet integrated evaluation tool on their respective devices. Next, the evaluation tool was used in accordance with the planned learning evaluation material flow. After all sessions were completed, students filled out a questionnaire containing statements to collect data related to the feasibility and effectiveness of the integrated system in supporting the measurement of their analytical skills in compiling Higher Order Thinking Skills questions.

Student Responses to Evaluation Tools

**Figure.15**. Graph of Student Responses to the Evaluation Tool

Referring to Figure 15, the deep learning-based Google Sheet integrated evaluation tool obtained an average feasibility of 90%, which places it in the highly feasible category at the trial use stage. Student responses indicate that the interface and workflow, from inputting Higher Order Thinking Skills items in Google Sheets to receiving analytical feedback from the model, are engaging, increase learning motivation, and facilitate understanding of the evaluation concept and how it works. This solution is considered appropriate for use in teaching evaluation material in the PPG program while also supporting the measurement of students' analytical skills in compiling HOTS questions.

### 4.1.5. Evaluation

Evaluation, as the final stage of the ADDIE model, is used to assess the achievement of the objectives of the evaluation tool that integrates deep learning-based Google Sheets into the Learning Evaluation material at PPG. The evaluation is carried out in layers at each phase, namely: First, Analysis, where researchers examine the results of mapping students' complexity and needs to ensure that the solution is relevant to the PPG context. The design is reviewed in terms of learning objective formulation, material development, questionnaire and test instrument design, and interaction flow storyboard from filling in Higher Order Thinking Skills items to result visualization, so that all elements support learning effectiveness. Development: expert validators assess the questionnaire's feasibility and identify shortcomings, such as typos, inaccurate illustrations, and the need to enrich the material to encourage higher-level thinking. Implementation stage: the questionnaire is filled out by practitioners and PPG students after using the evaluation tool. Practitioners gave positive responses regarding its effectiveness but suggested variations in application scenarios, while students reported an increase in understanding despite the guidebook's lack of detailed images that could be enlarged. This series of evaluations ensures that the development follows the ADDIE principles, resulting in effective teaching materials and evaluation mechanisms for Higher Order Thinking Skills that are tailored to student needs.

### 4.1.6. Effectiveness Analysis of Evaluation Tools

The effectiveness test was designed to assess the success of the deep learning-based Google Sheet integrated evaluation tool in improving the learning outcomes and analytical skills of PPG students in developing Higher Order Thinking Skills questions. The design used was a pretest-posttest with experimental and control groups. In the first meeting, all students underwent a pretest to map their initial abilities. In the second meeting, the experimental group learned using a Google Sheet-based workflow, including entering Higher Order Thinking Skills items, automatic processing by a deep learning model, and analytical feedback, while the control group followed conventional learning using printed teaching materials. In the third meeting, both groups took a posttest to test the improvement in understanding and learning achievement after the intervention, comparing the impact of the deep learning-based Google Sheet evaluation tool approach with the conventional method.

**Table.2**. Normality Test

| Treatment | Sig. |
|---|---|
| Experimental Pretest | 0.095 |
| Control Pretest | 0.136 |
| Experimental Posttest | 0.490 |
| Control Posttest | 0.198 |

The learning outcome data on the deep learning-based Google Sheet evaluation tool meet the normality assumption (Table 2). Using the Kolmogorov-Smirnov test, the experimental group showed pretest and posttest significance values of 0.095 and 0.490, respectively, while the control group showed pretest and posttest significance values of 0.136 and 0.195, respectively. All Sig. values were > 0.05, so the distributions of pretest and posttest data in both groups were deemed normal and suitable for further inferential analysis.

**Table.3**. Homogeneity Test

| | Levene Statistic | df1 | df2 | Sig. |
|---|---|---|---|---|
| Effectiveness of Evaluation Tools | 0.021 | 1 | 68 | 0.845 |

The results of the homogeneity test on the pretest and posttest data in both the experimental and control groups showed a significance value of 0.845. Based on statistical testing criteria, because the significance value of 0.845 was> 0.05, the data variances in both groups were deemed homogeneous (Table 2). With the prerequisites of homogeneity and normality of the data previously tested, statistical analysis can proceed using a parametric test (t-test) to test the hypothesis of the effectiveness of the Deep Learning-based Google Sheet evaluation tool on PPG students' ability to compile HOTS questions.

The results of the paired t-test in SPSS show a two-tailed p-value of 0.045 (< 0.05), so $H_0$ is rejected, and Ha is accepted. These findings confirm that evaluation tools that integrate deep learning-based Google Sheets are effective in improving learning outcomes and in supporting PPG students' analytical skills in compiling HOTS questions, compared with before the intervention. This integrated solution is worthy of implementation in educational evaluation learning in the Teacher Professional Education Program

### 4.2. Discussion

The use of Google Sheets and deep learning integration in this study not only serves as a variation in learning media but also plays a major role as a precision measurement tool for analytical skills. Technically, this system works by diagnosing the sentence structure in the questions to automatically determine their cognitive level, whether HOTS or LOTS. Direct feedback from AI trains students' logic to distinguish, organize, and compile complex question elements. Thus, the improvement of student competence is based on clear and measurable cognitive parameters, not merely an increase in motivation. In line with this, digital media can encourage active participation, increase enthusiasm, and positively impact learning outcomes (Muzammil & Izzah, 2025). The implementation of interactive formats, such as evaluation tools, significantly increases motivation and academic achievement, making them a relevant reference for strengthening the design of deep learning-based Google Sheets for Learning Evaluation materials in PPG (Banda & Nzabahimana, 2022).

The development of digital evaluation tools in this study is manifested as a deep learning-based Google Sheet evaluation tool for compiling Higher Order Thinking Skills questions, emphasizing the crucial role of web-based learning design. Microlearning through web pages helps students reflect on the knowledge they have learned while consolidating it in structured micro-learning units (Lawson et al., 2024), an approach that is in line with the needs of Learning Evaluation material. This involves understanding concepts in depth and then applying them practically to test items. In line with this, presenting material with digital tools makes the learning process more enjoyable, reduces boredom, and facilitates the understanding of complex concepts, a finding that is relevant for designing input flows in Google Sheets, automatic feedback from deep learning models, and user-friendly analytical displays (Shone et al., 2018).

In line with previous findings, evaluation tools that integrate deep learning-based Google Sheets as an extension of interactive approaches are not only effective in reducing boredom but also proven to significantly increase student motivation and learning achievement. Its visual-interactive design encourages active student participation by entering Higher Order Thinking Skills items into Google Sheets, receiving analytical feedback from the model, and reflecting on improvements iteratively. This more engaging and participatory learning experience aligns with the mainstream of web-based learning, which is recognized as innovative and effective, making it relevant for application in Learning Evaluation materials in PPG to strengthen students' analytical skills in developing Higher Order Thinking Skills questions.

### 5. Conclution

Based on the findings of the research, the design of an evaluation tool that integrates Google Sheets based on deep learning is in line with the results of the analysis of complexity and user needs in the Learning Evaluation material at PPG. Statistical testing using a paired t-test showed that the Sig. value was $0.000 < 0.05$, so Ha was accepted, and $H_0$ was rejected. This confirms the effectiveness of deep learning-based Google Sheet evaluation tools for students in the experimental class, with significantly greater gains than in the control class, as reflected in the larger difference between pretest and post-test scores. In the field, several obstacles arose, such as dependence on Google Sheets and connectivity, which could potentially affect input quality and synchronization stability; limited coverage of items and sub-materials, so that generalization to other Higher Order Thinking Skills domains has not been tested; insufficient transparency of model decisions on the dashboard to explain the reasons for assessment; and minimal adaptive formative evaluation and

collaboration features in Google Sheets. Future recommendations include expanding the material, adding a module explaining model decisions, enriching adaptive collaboration and evaluation features, strengthening data governance, and replicating across institutions and topics. This is expected to improve the accuracy, accountability, and effectiveness of this approach in measuring and strengthening the analytical abilities of PPG students in compiling Higher Order Thinking Skills questions.

## 6. References

Akdemir, E., & Barışçı, N. (2024). A review on deep learning applications with semantics. *Expert Systems with Applications*, *251*(1), 124029. https://doi.org/https://doi.org/10.1016/j.eswa.2024.124029

Aniek, O., Kusuma, S., Gede Ratnaya, I., & Candiasa, I. M. (2022). Evaluasi Pembelajaran Hybrid Pada Institusi Pendidikan Instiki Dengan Discrepancy Evaluation Model. *Lembaga Penjaminan Mutu*, *8*(2), 217-227. https://doi.org/10.25078/jpm.v8i02.1509

Ardiansyah, Risnita, & Jailani, M. S. (2023). Teknik Pengumpulan Data Dan Instrumen Penelitian Ilmiah Pendidikan Pada Pendekatan Kualitatif dan Kuantitatif. *Jurnal IHSAN : Jurnal Pendidikan Islam*, *1*(2), 1–9. https://doi.org/10.61104/ihsan.v1i2.5

Banda, H. J., & Nzabahimana, J. (2022). The Impact of Physics Education Technology (PhET) Interactive Simulation-Based Learning on Motivation and Academic Achievement Among Malawian Physics Students. *Journal of Science Education and Technology*, *32*, 127–141. https://doi.org/10.1007/s10956-022-10010-3

Bujang, M., Omar, E., & Baharum, N. (2018). A Review on Sample Size Determination for Cronbach's Alpha Test: A Simple Guide for Researchers. *The Malaysian Journal of Medical Sciences : MJMS*, *25*(6), 85–99. https://doi.org/10.21315/mjms2018.25.6.9

Chen, G., & Han, Z. (2025). The rise of deep learning: AI and engineering applications under the spotlight of the 2024 Nobel prize. *Intelligent Geoengineering*, *2*(1), 14–21. https://doi.org/10.1016/j.ige.2025.03.002

Dagal, I., Harrison, A., Ibrahim, A. W., & Mbasso, W. F. (2025). Comprehensive evaluation of data preprocessing and visualization techniques for enhanced classification and sampling. *Cluster Computing*, *28*(7), 476. https://doi.org/10.1007/s10586-025-05512-9

Dale, E. (1969). *Audiovisual methods in teaching*. New York: Holt, Rinehart and Winston Inc. The Dryden Press.

DeBell, T., Goertzen, L., Larson, L., Selbie, W., Selker, J., & Udell, C. (2019). OPEnS hub: Real-time data logging, connecting field sensors to google sheets. *Frontiers in Earth Science*, *7*, 1-6. https://doi.org/10.3389/feart.2019.00137

Dos Santos, L. R., Behr, A., & Duarte, G. R. (2025). Recognizing accounting as a STEM discipline through professional skills in accounting information systems☆. *Journal of Accounting Education*, *72*. https://doi.org/10.1016/j.jaccedu.2025.100970

Gardner, J., O'Leary, M., & Yuan, L. (2021). Artificial intelligence in educational assessment: 'Breakthrough? Or buncombe and ballyhoo?' *Journal of Computer Assisted Learning*, *37*(5), 1207–1216. https://doi.org/10.1111/jcal.12577

Gkintoni, E., Antonopoulou, H., Sortwell, A., & Halkiopoulos, C. (2025). Challenging Cognitive Load Theory: The Role of Educational Neuroscience and Artificial Intelligence in Redefining Learning Efficacy. In *Brain Sciences, 15*(2), Multidisciplinary Digital Publishing Institute (MDPI). https://doi.org/10.3390/brainsci15020203

Huang, Z., Fu, X., & Zhao, J. (2025). Research on AIGC-Integrated Design Education for Sustainable Teaching: An Empirical Analysis Based on the TAM and TPACK Models. *Sustainability (Switzerland)*, *17*(12). https://doi.org/10.3390/su17125497

Kacmaz, G., & Dubé, A. K. (2025). Measuring teacher scaffolding in game-based learning: Emotional and responsibility scaffolds lead while resources and previous game experience shape practices. *Computers and Education*, *237*. https://doi.org/10.1016/j.compedu.2025.105393

Kasanah, M., & Pratama, P. A. (2024). Taksonomi Tujuan Pendidikan dan Evaluasi Hasil Belajar. *Pedagogik: Jurnal Pendidikan dan Riset*, *2*(2), 146–162. https://ejournal.edutechjaya.com/index.php/pedagogik/article/view/593

Khoiriyah, Z., & Shaleh, S. (2025). Solusi Alternatif Atas Problematika dalam Mengimplementasikan Penilaian Higher Order Thingking Skills (HOTS). *Al-Madrasah Jurnal Pendidikan Madrasah Ibtidaiyah*, *9*(2), 656. https://doi.org/10.35931/am.v9i2.4855

Kunicki, Z. J., Zambrotta, N. S., Tate, M. C., Surrusco, A. R., Risi, M. M., & Harlow, L. L. (2019). Keep Your Stats in the Cloud! Evaluating the Use of Google Sheets to Teach Quantitative Methods. *Journal of Statistics Education*, *27*(3), 188–197. https://doi.org/10.1080/10691898.2019.1665485

Lawson, A. A., Ramage, K., Gibson, M. E. S., King, C., Imbo-Nloga, C. C., & French, A. V. (2024). Pediatric and Adolescent Gynecology WebEd: A Brief Report of an Underutilized Online Learning Module. *Journal of Pediatric and Adolescent Gynecology*, *37*(3), 371–374. https://doi.org/https://doi.org/10.1016/j.jpag.2024.01.166

Li, X., & Zhang, Y. (2025). Security application of intrusion detection model based on deep learning in english online education. *Alexandria Engineering Journal*, *124*, 582–590. https://doi.org/10.1016/j.aej.2025.03.051

Ligasan, J. M., & Doysabas, J. C. (2024). Enhancing news writing skills of grade 4 students through graphic organizer strategy. *Davao Research Journal*, *15*(4), 83-94. https://doi.org/10.59120/drj.v15i4.275

Ma, Y., Sun, X., & Ma, A. (2025). On the assessment and reliability of political and ideological education in colleges using deep learning methods. *Alexandria Engineering Journal*, *119*, 511–517. https://doi.org/10.1016/j.aej.2025.01.114

Maulana, F. F., Deffinika, I., Jazilah, F., & Pokhrel, L. B. (2023). Virtual Learning Media Exploration 3D View of the Konto River Watershed to Support 21st Century Learning.

*Jambura Geo Education Journal*, *4*(2), 133–146. https://doi.org/10.34312/jgej.v4i2.21399

Maxnun, L., Kristiani, K., & Sulistyaningrum, C. D. (2024). Development of hots-based cognitive assessment instruments: ADDIE model. *Journal of Education and Learning (EduLearn)*, *18*(2), 489–498. https://doi.org/10.11591/edulearn.v18i2.21079

Mir, S., Shah, S. A., & Bhat, M. S. (2024). Evaluation of disaster preparedness in educational infrastructure of remote geographies using Analytical Hierarchy process. *International Journal of Disaster Risk Reduction*, *108*. https://doi.org/10.1016/j.ijdrr.2024.104514

Misnawati, M., Junari, J., Teibang, D., Ilham, I., & Luthfiyah, L. (2025). Evaluasi Hasil Asesmen Melalui Pemberian Umpan Balik dalam Tes Formatif sebagai Tolak Ukur Hasil Belajar Siswa. *JIIP-Jurnal Ilmiah Ilmu Pendidikan*, *8*(2), 2236–2242. https://doi.org/10.54371/jiip.v8i2.6836

Muzammil, M., & Izzah, S. A. (2025). Penerapan Pembelajaran Berdiferensiasi Dengan Media Digital Untuk Meningkatkan Hasil Belajar Siswa. *Edukasiana: Jurnal Inovasi Pendidikan*, *4*(3), 1237–1246. https://doi.org/10.56916/ejip.v4i3.1407

Nasbi, I. (2017). MANAJEMEN KURIKULUM: Sebuah Kajian Teoritis. *Idaarah: Jurnal Manajemen Pendidikan*, *1*(2), 318–330. https://doi.org/10.24252/idaarah.v1i2.4274

Nasser, M., Arshad, N. I., Ali, A., Alhussian, H., Saeed, F., Da'u, A., & Nafea, I. (2025). A systematic review of multimodal fake news detection on social media using deep learning models. In *Results in Engineering, 26*, Elsevier B.V. https://doi.org/10.1016/j.rineng.2025.104752

Nasution, A. fattah. (2023). *Metode Penelitian kualitatif*. Harva creative. https://harfacreative.com/

Parra, F., Jacobs, A., & Trevino, L. L. (2021). Shippy Express: Augmenting accounting education with Google Sheets. *Journal of Accounting Education*, *56*. https://doi.org/10.1016/j.jaccedu.2021.100740

Rebman, C. M., Booker, Q. E., Wimmer, H., Levkoff, S., Mcmurtrey, M., & Powell, L. M. (2023). An Industry Survey of Analytics Spreadsheet Tools Adoption: Microsoft Excel vs Google Sheets. In *Information Systems Education Journal (ISEDJ), 21*(5). https://isedj.org/;https://iscap.infohttps://isedj.org/;https://iscap.info

Sandika, N. A., Basit, A., Putra, R. J., Nufrizal, A., Putra, B. A., Gunawan, R. G., Ameliya, D., & Utami, R. J. (2025). *Deep Learning dalam Pendidikan Tinggi: Strategi, Metode, dan Implementasi*. Pradina Pustaka.

Saputra, D. A., & Nofrion. (2022). Pengembangan LMS Berbasis Google Sites untuk Mengembangkan HOTS Siswa pada Mata Pelajaran Geografi. *Jurnal Pendidikan Ilmu Pengetahuan Sosial Indonesia*, *7*(3), 206–220. https://journal.stkipsingkawang.ac.id/index.php/JurnalPIPSI/article/view/3147

Setiawan, Z., Pustikayasa, I. M., Jayanegara, I. N., Setiawan, I. N. A. F., Putra, I. N. A. S., Yasa, I. W. A. P., Asry, W., Arsana, I. N. A., Chaniago, G. G., Wibowo, S. E., & others.

(2023). *Pendidikan Multimedia : Konsep dan Aplikasi pada era revolusi industri 4.0 menuju society 5.0.* PT. Sonpedia Publishing Indonesia. https://books.google.co.id/books?id=7YjhEAAAQBAJ

Shone, N., Ngoc, T. N., Phai, V. D., & Shi, Q. (2018). A deep learning approach to network intrusion detection. *IEEE Transactions on Emerging Topics in Computational Intelligence*, *2*(1), 41–50. https://doi.org/10.1109/TETCI.2017.2772792

Sugrah, N. U. (2020). Implementasi teori belajar konstruktivisme dalam pembelajaran sains. *Humanika*, *19*(2), 121–138. https://doi.org/10.21831/hum.v19i2.29274

Sumartini, S., Harahap, K. S., & Sthevany, S. (2020). Kajian Pengendalian Mutu Produk Tuna Loin Precooked Frozen Menggunakan Metode Skala Likert Di Perusahaan Pembekuan Tuna. *Aurelia Journal*, *2*(1), 29. https://doi.org/10.15578/aj.v2i1.9392

Sumy, D. F., Drakes, O. O., McBride, S. K., & Jenkins, M. R. (2023). Social vulnerability and geographic access barriers to earthquake early warning education in museums and other free choice learning environments. *International Journal of Disaster Risk Reduction*, *97*. https://doi.org/10.1016/j.ijdrr.2023.104011

Surahman, E., Soepriyanto, Y., Wedi, A., & Ulfa, S. (2024). *Online Assessment: Panduan Guru Sekolah Dasar*. Academia Publication.

Vanhees, C., Nijlunsing, J., Muijs, D., Crato, N., Wils, M., Wiliam, D., Surma, T., & Kirschner, P. A. (2025). The role of knowledge-rich curricula in promoting deep thinking and complex skill acquisition. *Learning and Individual Differences*, *121*. https://doi.org/10.1016/j.lindif.2025.102729

Wahyudi, S., Arifin, T., Mutia, C., & Aji, H. (2025). Integrating Publuu technology into interactive teaching materials as. In *Jurnal Inovasi Teknologi Pendidikan, 21*(1). https://doi.org/10.21831/jitp.v12i1.78293

Yudistira, E. N., & Kom, S. (2024). *Deep Learning: Teori, Contoh Perhitungan, dan Implementasi*. Deepublish.