

Identification of Book Cover Titles Using the Natural Language Processing (NLP) Method

Rianti Afifah Gultom^{1*}, Munjiat Setiani Asih², Ade Zulkarnain Hasibuan³

^{1,2} Faculty of Engineering and Computer Science, Department of Informatics Engineering, University of Harapan Medan, Medan, Indonesia

³ Faculty of Science and Technology, Department of Informatics, University of Samudra, Langsa, Indonesia
Author(s) Email: ¹rafifahgultom@gmail.com, ²munjiat.stth@gmail.com, ³adezulhsb@unsam.ac.id

ARTICLE INFO

Article history:

Received November 11, 2025

Revised November 11, 2025

Accepted November 11, 2025

Publish November 30, 2025

ABSTRACT

In the digital era, the identification of book titles on covers has become a crucial requirement in digital library management, archiving systems, and book e-commerce platforms. The main challenges lie in the limitations of manual methods and traditional pattern-matching techniques, which are inefficient, as well as in the complexity of processing the Indonesian language, which exhibits diverse morphological variations and syntactic structures. To address these issues, this study proposes the integration of Optical Character Recognition (OCR) with the Natural Language Processing (NLP) method. OCR is utilized to extract textual information from book cover images, while NLP is applied to recognize and classify the extracted text to identify the main book title. The implementation results demonstrate that this approach significantly improves title identification accuracy compared to traditional methods, particularly through the application of Named Entity Recognition (NER) techniques and modern NLP models such as BERT and LSTM. The developed system proves effective in accelerating the book digitalization process, enhancing information management efficiency, and contributing to the advancement of Indonesian language processing technology.

Keywords:

Optical Character Recognition (OCR), Natural Language Processing (NLP), Book Title Identification, Named Entity Recognition (NER), Indonesian Language.

Corresponding Author:

Rianti Afifah Gultom,

Faculty of Engineering and Computer Science, Department of Informatics Engineering, Universitas Harapan Medan, Medan, Indonesia

HM. Joni Street No. 70C

Email: rafifahgultom@gmail.com



1. INTRODUCTION

In the digital era, the growing volume of digital documents necessitates automated processing, including the identification of book titles on book covers. One promising approach to address this problem is Natural Language Processing (NLP), a technique that enables machines to analyze, understand, and extract information from human language text. Book title identification is a crucial component in digital library management systems, archiving systems, and even book e-commerce applications. The main challenge addressed in this study lies in the automatic identification of book titles on book covers. In digital library management, archiving systems, and book e-commerce applications, manual identification or traditional pattern-matching methods have become inefficient. This challenge becomes even more complex when dealing with natural languages such as Indonesian, which exhibits unique morphological variations and structural characteristics. Therefore, a more advanced approach, such as Natural Language Processing (NLP), is required to improve the accuracy and efficiency of automatic book title recognition.

The use of NLP in book title identification can greatly facilitate the classification, retrieval, and management of large-scale book data. NLP can be employed to extract key information such as book titles by utilizing techniques like Named Entity Recognition (NER) and Text Classification. This technology enables systems to recognize important entities within text such as titles, authors, and publishers thereby making information management more efficient and structured. Moreover, the application of NLP models in complex natural languages such as Indonesian has shown increasingly promising results across various information-related applications.

They have shown promising results in various information applications. In the context of text processing in the Indonesian language, the development of NLP-based systems for book title identification on covers has not yet been extensively explored. A previous study conducted by Fajri Koto et al. (2020) in their work “IndoLEM and IndoBERT: A Benchmark Dataset and Pre-trained Language Model for Indonesian NLP” demonstrated that Indonesian-specific pre-trained models can achieve state-of-the-art performance across various morpho-syntactic, semantic, and discourse tasks. Nevertheless, several NLP models that have been developed, such as BERT and LSTM, exhibit great potential for processing Indonesian-language text.

Furthermore, the identification of book titles on covers serves as an essential initial step in the broader process of book digitalization. Digital library systems that rely on automatic book title recognition demonstrate greater speed and efficiency in archive management compared to manual systems. This approach also highlights how the application of NLP can enhance efficiency in book retrieval and archiving, as well as accelerate the digitalization process of older books that have not yet been integrated into digital systems.

However, the application of Natural Language Processing (NLP) in identifying book titles on Indonesian-language covers still faces several challenges. The Indonesian language possesses distinct morphological and syntactic characteristics compared to other languages, necessitating adaptations in the NLP models being utilized. Although NLP models such as Bidirectional Encoder Representations from Transformers (BERT) and Long Short-Term Memory (LSTM) have demonstrated promising results in text processing, their implementation for the Indonesian language—particularly in the context of book title recognition—has not yet been extensively explored.

Overall, this study aims to develop an NLP-based method for identifying titles on Indonesian-language book covers. By employing various well-established NLP techniques, the proposed system is expected to enhance the efficiency of book information management and address the challenges of natural language processing in the context of the Indonesian language. The effective use of NLP in book title recognition on covers can also help accelerate the integration of book information into larger library systems or digital platforms [1].

Therefore, this study aims to develop an NLP-based method capable of automatically identifying book titles on Indonesian-language book covers. By leveraging NLP techniques such as Named Entity Recognition (NER) and Text Classification, the proposed system is expected to address challenges in processing Indonesian text and enhance the efficiency of book information management. The implementation of this system is also expected to support a broader book digitalization process, benefiting both digital library systems and book e-commerce platforms.

2. RESEARCH METHODOLOGY

2.1 Title Identification

Title identification is the process of recognizing, extracting, and classifying text that serves as the main name of a document or medium, such as a book. In the context of digital information systems and Natural Language Processing (NLP) [2], title identification aims to automatically detect the portion of text that represents the title from data sources such as book covers or textual metadata [3].

Book titles often exhibit distinctive characteristics, such as the use of capital letters, placement at the top or center of the cover, and word choices that reflect the book’s main topic. However, in practice, challenges arise due to variations in cover design, layout, font, and linguistic structure particularly in the Indonesian language, which is morphologically rich. Therefore, an NLP-based approach is required to accurately recognize the “title” entity using techniques such as Named Entity Recognition (NER), Text Classification, or deep learning models such as BERT and LSTM.

Several approaches also incorporate visual pre-processing methods, such as Optical Character Recognition (OCR), to extract text from cover images, which is subsequently subjected to linguistic analysis. The integration of OCR and Named Entity Recognition (NER) can enhance the system’s accuracy in identifying valid title text.

2.2 Machine Learning

Machine learning is an analytical process designed to discover data patterns and relationships among data variables [4]. One of its primary features is the ability to analyze complex non-linear relationships, as well as to accommodate highly intricate input variables. Numerous machine learning models can be adapted to analyze data through classification, clustering, and association rule mining, depending on the suitability of the collected data and the objectives of the analytical process [5]. The main advantage of machine learning is that once a sufficiently large dataset has been gathered, the algorithm can learn how to handle the available data, after which it is able to operate automatically [6].

Machine learning is useful for monitoring and analyzing learning processes in schools, predicting student performance by providing necessary academic support, academic guidance, and mentorship, evaluating the efficiency and effectiveness of teaching methods, as well as delivering meaningful feedback to both teachers and students. It can also modify learning environments to benefit students based on information or data related to their background issues and academic progress throughout a semester. This demonstrates that machine learning can be effectively utilized to predict student development across various educational levels [7].

Machine learning offers significant advantages in the modern era compared to traditional forms of statistical analysis, as it emphasizes predictive performance over proven theoretical properties and a priori population assumptions. Machine learning is employed to achieve this goal. Its techniques are used to discover models or data patterns, which in turn assist in decision-making processes. The ability to predict student performance is particularly crucial in today's educational systems. However, it remains unclear which machine learning models are most effective in predicting student performance and which are best suited for improving learning outcomes [8].

There are various data mining methods commonly used to extract hidden information from large volumes of data. Machine learning models such as decision trees, neural networks, Bayesian classifiers, nearest neighbors, support vector machines, random forests, logistic regression, linear discriminant analysis, multiple regression, and self-organizing maps are among the most widely applied. Machine learning is a computational approach that enables systems to learn patterns from data and make predictions or decisions without being explicitly programmed [9].

2.3 Natural Language Processing (NLP)

Natural Language Processing (NLP) is a rapidly evolving field within computer science and linguistics. NLP is concerned with the understanding, interpretation, and generation of human language by computer systems [10]. One of the foundational theories underlying NLP is the Theory of Language Processing developed by Noam Chomsky in the 1950s [11]. This theory posits that human language possesses an inherent structural regularity and that language comprehension depends on specific syntactic and semantic rules [2].

In addition, the Theory of Knowledge Representation also serves as an important foundation in NLP. This theory posits that the meanings of words and phrases in a language represent real-world knowledge. For instance, the concept of word vector representations (word embeddings) in NLP is based on the idea that words frequently appearing together within the same context tend to share similar meanings [12].

In 2013, Mikolov and his colleagues developed the Word2Vec technique, which became one of the major breakthroughs in the development of word representations in NLP. This technique employs an artificial neural network model to generate vector representations of words based on their co-occurrence within similar contexts. It enables machines to better understand the relationships between words and their contextual meanings in text [13].

In recent years, the adoption of transformer-based language models such as BERT (Bidirectional Encoder Representations from Transformers), developed by Google in 2018, has revolutionized the field of NLP. This model is capable of generating more refined word representations by considering both the preceding and succeeding contexts within a sentence, allowing machines to comprehend text at a higher level of understanding [14]. After the text is successfully extracted from book cover images through the OCR process, the next stage involves the application of Natural Language Processing (NLP) to process, analyze, and identify relevant text components such as the main title, author name, or other pertinent information. This process is carried out through a series of steps as follows:

1. Text Extraction (OCR)
Extracting text from book cover images using Optical Character Recognition.
2. Text Preprocessing
 - a. Case Folding: Converting all letters to lowercase.
 - b. Cleaning: Removing irrelevant symbols or punctuation marks.
 - c. Tokenization: Splitting the text into individual words.
 - d. Stopword Removal: Eliminating common words that do not carry analytical meaning.
 - e. Stemming/Lemmatization: Reducing words to their base or root form.
3. Linguistic Analysis
 - a. POS Tagging: Assigning part-of-speech labels (nouns, verbs, adjectives, etc.).
 - b. Named Entity Recognition (NER): Detecting important entities such as titles, authors, or publication years.
4. Title Classification
Determining the text segment most likely to represent the main title based on linguistic patterns and contextual cues.
5. Output
Displaying the identified title along with the necessary description or metadata.

2.4 OCR (Optical Character Recognition)

Optical Character Recognition (OCR) is a technology used to convert text contained in images—such as text on book covers, receipts, or other printed documents—into digital text that can be read and processed by a computer. In the context of the study “Identification of Book Cover Titles Using the Natural Language Processing (NLP) Method,” OCR serves as a crucial initial stage, as it enables the extraction of text from book covers prior to further processing using NLP methods to identify the actual title [15].

In the context of book cover title identification, OCR is utilized to read text from cover images, which is then further processed using NLP techniques such as Named Entity Recognition (NER) or Text Classification to recognize the text segment representing the main title. The effectiveness of OCR is highly influenced by image quality, font type, contrast, and the presence of noise or distortions on the cover [16].

Modern OCR technologies, such as Tesseract OCR, EasyOCR, and Google Vision OCR, have been widely applied in various domains, including documentation systems, digital archives, and digital libraries. OCR can also be integrated with deep learning techniques to enhance accuracy under complex image conditions, for example by employing Convolutional Neural Networks (CNN) for character segmentation. The OCR process generally involves the following steps:

1. Image Input: Capture the book cover image.
2. Preprocessing: Convert the image to grayscale, remove noise, and correct alignment.
3. Text Extraction (OCR): Use Tesseract or EasyOCR to recognize text from the image.
4. Post-processing: Correct and clean the extracted OCR text.
5. NLP Analysis: Perform tokenization and Named Entity Recognition to detect the book title.
6. Output: Display the detected text as the book title.

3. RESULT AND DISCUSSION

3.1 Analysis

In the development of a book-cover title identification system based on Natural Language Processing (NLP), several issues must be analyzed to ensure optimal system performance. A primary challenge is the heterogeneity of cover design formats: each book may present distinct layouts, font sizes, and visual elements. Such variability complicates the automation process for accurately determining which text segment constitutes the principal title.

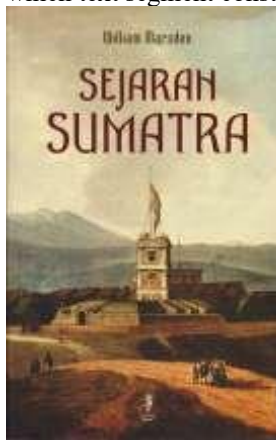


Figure 1. Book Cover

The next issue concerns the accuracy of text extraction results using OCR. Image quality, resolution, and the presence of non-text elements such as graphics or symbols can significantly affect the success of the text extraction process, which serves as the primary input for subsequent NLP processing.

Moreover, there are challenges in developing an NLP model capable of distinguishing book titles from other textual elements such as author names, publishers, or promotional slogans. This process requires a learning-based or linguistically rule-driven approach to accurately identify the linguistic structures commonly used in book title composition.

Another crucial issue that must be taken into account involves data management and system security, particularly in ensuring that only authorized users can access and manage the results of the identification process. By systematically identifying these challenges, system development can be directed toward appropriate solutions aimed at enhancing the accuracy and efficiency of the book cover title identification process.

3.1.1 Book Cover Analysis

Based on the identification results from this book cover, the system should be able to conclude that the text "SEJARAH SUMATRA" represents the main title of the book, as indicated by its larger font size, central placement, and linguistic context. Accurate identification at this stage is crucial, as any error in detecting the title may affect the overall quality and integrity of the book's metadata.

Table 1. Book Cover Identification

Element	Description
Book Title	<i>SEJARAH SUMATRA</i> — Located in the center, large font size, strong contrast color
Author Name	William Marsden — Positioned at the top, smaller font size
Visual Elements	Illustration of fortress, flag, and historical icons — does not contain text
Background Color	Natural and historical tones — supports classical visual aesthetics
Identification Focus	The book title is determined based on its position, font size, and linguistic context

Table 2. Stages of OCR and NLP Processes

Stage	Process Description	Recommended Tool
1 Image Preprocessing	Adjust contrast, remove noise, and convert to grayscale	OpenCV
2 Text Extraction (OCR)	Convert text from image into digital text	Tesseract OCR
3 Tokenization	Split OCR text results into tokens/words for further analysis	spaCy, NLTK
4 Text Cleaning	Remove special characters, stopwords, and OCR noise	NLTK, regex
5 Named Entity Recognition (NER)	Identify named entities such as title, author, or location	spaCy, NLTK
6 Text Classification (Optional)	Use NLP model or ML classifier to determine which text represents the title	Custom NLP rule-based or ML classifier
7 Title Output	Display text with the highest probability of being the title	Web Frontend / System Output

Based on the book cover analysis, the system can determine that the text "SEJARAH SUMATRA" represents the main title of the book. This conclusion is supported by several visual and linguistic indicators, including the text's central position on the cover, the larger font size compared to other elements, and the use of contrasting colors that make it stand out among other visual components. In addition, supporting information such as the author's name positioned at the top in smaller font size and visual illustrations without text further strengthens the title identification. The primary focus at this stage is to ensure accurate title recognition, as errors in early-stage identification may affect the overall accuracy of the book's metadata.

The title identification process utilizes a combination of Optical Character Recognition (OCR) and Natural Language Processing (NLP) techniques. The stages begin with image pre-processing to enhance text quality, followed by text extraction using Tesseract OCR, tokenization, and text cleaning to remove noise. Subsequently, the system applies Named Entity Recognition (NER) to identify key entities such as the book title and author name, along with an optional text classification step to strengthen title detection. The final result is an output text with the highest probability of being the book's title, which is then displayed through the system or a web-based interface [17].

At this stage, a discussion is conducted on the process of title identification on book covers based on the system workflow and the implementation results of Natural Language Processing (NLP) techniques. Based on the example of the book cover "Sejarah Sumatra" by William Marsden, it can be observed that the system must be capable of distinguishing text elements such as the title and the author's name by analyzing both visual features (size, position, color) and linguistic characteristics of the recognized text.

After the book cover image is processed using Optical Character Recognition (OCR), the extracted text is then analyzed using Natural Language Processing (NLP) to recognize sentence structures and predict which segment represents the main title. In this case, the text "SEJARAH SUMATRA" is successfully identified as the title because it meets several indicators, such as large font size, strategic central placement, and a phrase structure commonly used for book titles.

This discussion demonstrates that the integration of Optical Character Recognition (OCR) and Natural Language Processing (NLP) can produce a fairly accurate title identification process. However, its success is highly dependent on the image quality and the presence of other textual elements that may introduce ambiguity. Therefore, this process needs to be refined through additional approaches such as machine learning-based text pattern recognition or weighting specific visual features to improve accuracy in determining the book's main title.

3.2 System Implementation

The interface display represents the stage where the system is ready to be operated under real-world conditions in accordance with the results of prior analysis and design. Through this stage, it can be determined whether the developed system successfully meets its intended objectives.

The application is equipped with an interface designed to facilitate user interaction with the system. The primary function of the interface is to receive input from the user and display the output produced by the identification process. In the Book Cover Title Identification application using the NLP method, the interface consists of a login form, a book cover upload form, and a title identification form that processes the input using NLP algorithms.

The interface display represents the stage where the system is ready to be operated under real-world conditions in accordance with the results of prior analysis and design. Through this stage, it can be determined whether the developed system successfully meets its intended objectives. The application is equipped with an interface designed to facilitate user interaction with the system. The primary function of the interface is to receive input from the user and display the output produced by the identification process. In the Book Cover Title Identification application using the NLP method, the interface consists of a login form, a book cover upload form, and a title identification form that processes the input using NLP algorithms.

1. Login Form

The login form functions as the entry gateway to the system, ensuring security by preventing unauthorized user access before entering the main menu. This form consists of input fields for the username and password, along with a button to submit the login data. The required attribute is applied to both input fields to ensure that no data is left empty before submission. Once the user fills out and submits the login form, the system processes the data to verify user authenticity.

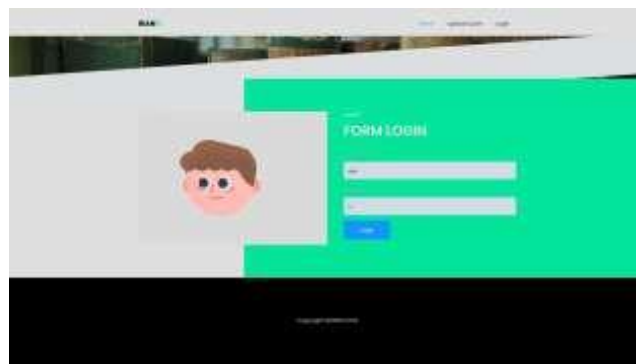


Figure 2. Login Form

2. Main Menu Form

The main menu form serves as a central interface that connects to other components of the system, including the book data form and the NLP algorithm processing form. Through this form, users can easily navigate between different system features such as uploading book covers, managing book data, and executing the title identification process. Below is the display of the main menu form:



Figure 3. Main Menu

The Book Data Form is utilized by the administrator to manage and store data within the database. This form provides functionalities to add new book records, manage related variables, and execute the NLP algorithm process. Through this form, administrators can organize book data systematically, ensuring that the system operates efficiently and optimally in identifying book titles from cover images.

- [3] T. D. W. Negara, "Analisis Desain Cover Buku Baca Anak Usia Dini Karya Gibran Maulana," 2021.
- [4] A. Wibowo and A. R. Isnain, "Implementasi Algoritma Machine Learning untuk Klasifikasi Suara Lingkungan," *MALCOM Indones. J. Mach. Learn. Comput. Sci.*, vol. 5, no. 2, pp. 616–625, 2025.
- [5] A. Nugroho, R. Setiawan, A. Haris, and Beny, "Deteksi Bahasa Isyarat Bisindo Menggunakan Metode Machine Learning," *Process. J. Ilm. Sist. Informasi, Teknol. Inf. dan Sist. Komput.*, vol. 18, no. 2, pp. 152–158, 2023.
- [6] G. Rizaldy, M. S. Asih, and Y. F. A. Lubis, "Implementasi Jaringan Saraf Tiruan Untuk Segmentasi Dan Klasifikasi Daun Jambu Menggunakan Metode PCA Dan K-NN," *J. Ilmu Komput. Dan Teknol.*, vol. 1, no. 2, pp. 124–137, 2025.
- [7] H. Santoso and T. H. Rochadiani, "Pelatihan Machine Learning Menggunakan Bahasa Pemrograman Python Bagi Karyawan PT. Yokogawa Indonesia," *J. ABDINUS J. Pengabd. Nusant.*, vol. 6, no. 2, pp. 349–356, 2022.
- [8] M. F. Naufal and S. F. Kusuma, "Analisis Perbandingan Algoritma Machine Learning Dan Deep Learning Untuk Klasifikasi Citra Sistem Isyarat Bahasa Indonesia (SIBI)," *J. Teknol. Inf. dan Ilmu Komput.*, vol. 10, no. 4, pp. 873–881, 2023.
- [9] M. S. Asih, "Pengenalan Huruf Pada Citra Digital Menggunakan Algoritma Template Matching," *Jur. Tek. Inform. Sekol. Tinggi Tek. Harapan Medan*, vol. 1, no. 1, pp. 1–5, 2017.
- [10] I. Huda, "Implementasi Natural Language Processing(NLP) Untuk Aplikasi Pencarian Lokasi," *J. Nas. Teknol. Terap.*, vol. 3, no. 2, pp. 15–28, 2019.
- [11] A. Puspitasari, A. N. Paradhita, Y. W. Tineka, V. Sulistyowati, N. K. S. Noriska, and Haryanto, "Natural Language Processing (NLP) Technology for Chatbot Website," *J. Penelit. Pendidik. IPA*, vol. 10, no. 1, pp. 319–324, 2024.
- [12] I. Abdurrohman and A. Rahman, "Penerapan Natural Language Processing Untuk Analisis Sentimen Terhadap Kebijakan Pemerintah," *J. Kebangs. RI*, vol. 1, no. 1, pp. 55–60, 2024.
- [13] S. Rahmadani and E. Tasrif, "Implementasi Natural Language Processing (NLP) pada Layanan Pengelolaan Surat Kantor Camat Bukit Barisan," *URNAL Tek. Komput. DAN Inform.*, vol. 4, no. 1, pp. 19–24, 2024.
- [14] M. Amien, "Sejarah dan Perkembangan Teknik Natural Language Processing (NLP) Bahasa Indonesia: Tinjauan tentang sejarah, perkembangan teknologi, dan aplikasi NLP dalam bahasa Indonesia," *ELANG J. Interdiscip. Res.*, vol. 1, no. 1, pp. 99–105, 2023.
- [15] I. N. T. Lestari and D. I. Mulyana, "Implementation Of OCR (Optical Character Recognition) Using Tesseract In Detecting Character In Quotes Text Images," *J. Appl. Eng. Technol. Sci.*, vol. 4, no. 1, pp. 58–63, 2022.
- [16] A. Marshanda, B. Harijanto, and C. Rahmad, "Implementasi Optical Character Recognition (OCR) Untuk Meningkatkan Akurasi Dan Kecepatan Input Data Di Posyandu," *JIP (Jurnal Inform. Polinema)*, vol. 11, no. 1, pp. 45–50, 2024.
- [17] I. Hafaz, N. Wulan, and M. S. Asih, "Implementasi Algoritma Aes Dan Elgamal Untuk Enkripsi Dan Dekripsi File Data Video Mp4 Berbasis Web," *JATI (Jurnal Mhs. Tek. Inform.)*, vol. 9, no. 5, pp. 7964–7971, 2025.