15th August 2025. Vol. 103. No. 15 © Little Lion Scientific



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

# DEEP LEARNING-DRIVEN WEB INFORMATION EXTRACTION WITH CNN AND LSTM MODELS

B. BHAVANI<sup>1</sup>, Dr. D.HARITHA<sup>2</sup>

<sup>1</sup>Scholar, Department of CSE, Koneru Lakshmaiah Education Foundation (Deemed to be University), Vaddeswaram, AP, India.

<sup>2</sup>Professor, Department of CSE, Koneru Lakshmaiah Education Foundation (Deemed to be University), Vaddeswaram, AP, India.

Email: rbbhavanib2606@gmail.com<sup>1</sup>, haritha donavalli@kluniversity.in<sup>2</sup>

#### **ABSTRACT**

The high rate at which the amount of web data is growing has posed immense opportunities as well as challenges to the data extraction methods. Traditional web scraping solutions are easily outmatched by the dynamism and heterogeneity of web content leading to frequent failure and inefficiency. The article questions how these weaknesses could be mitigated by using adaptive deep learning models to extract web data by giving a robust and scalable solution to the problem. We propose an adaptive deep learning model learning and generalizing over diverse web structures and types of content using neural networks. It is a dynamic framework and can fit in changing web environments by utilizing domain adaptation and transfer learning to ensure consistency and accuracy in data extraction. We have determined through considerable experimentation that our adaptive architecture is considerably faster than the established scraping approach, with accuracy, structural change resilience and a reduction in manually configured dependency being noteworthy enhancements. These results illustrate the effectiveness of adaptive deep learning over web data extraction and lead to the path of much smarter and automated web scraping systems.

**Keywords:** Web Data, Conventional Methods, Deep Learning, LSTM, Web-Scraping

#### 1. INTRODUCTION

The exponential growth of the Internet has resulted in an unparalleled surge in the accessibility of information. The extensive collection of data, encompassing several fields like e-commerce, social media, scientific research, and others, offers significant prospects for making decisions and conducting analysis based on data. Nevertheless, the task of retrieving valuable and organized information from the internet continues to be a major obstacle because of the ever-changing, diverse, and frequently disorganized nature of online content. Conventional techniques for web scraping, which depend on predetermined rules and scripts, often become useless due to alterations in webpage layouts, structures, and technology. This requires regular manual upgrades and maintenance, which are both time-consuming and resource intensive. In order to overcome these constraints, current advancements in deep learning provide encouraging remedies. Deep learning models, especially those utilizing neural networks, have exhibited exceptional

achievements in many tasks related to pattern recognition, natural language processing, and computer vision. Due to their capacity to acquire knowledge from extensive datasets and apply it to novel, unfamiliar instances, they are highly suitable for the intricate undertaking of extracting data from the web. Nevertheless, employing deep learning directly for web scraping presents its own distinct set of obstacles, such as the requirement for extensive quantities of labelled training data and the issue of adjusting to always changing web material.

Figure 1 depicts the workflow of a web application interface, showcasing the interactions between users, administrators, and a server. The web browser functions as the client, while the server manages incoming requests and interactions. Users have the ability to access and see information about products, orders, and payments, whereas administrators possess the authority to oversee and control data. The login interface facilitates the process of verifying the identity and granting access privileges, while the database is responsible for storing and manipulating data associated with the online application.

© Little Lion Scientific



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

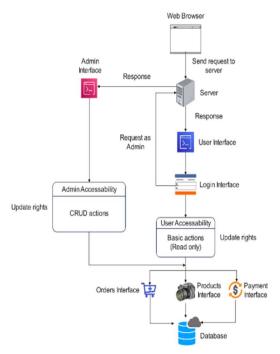


Figure 1: Existing Architecture Of Web Data Extraction

The current architecture presents various challenges, such as scalability limitations, a single point of failure, restricted ability to handle dynamic content, vulnerabilities arising from security communication between user and admin interfaces, inflexibility in role-based access control, reliance on CRUD operations, absence of advanced data processing or machine learning capabilities. inefficient utilization of distributed computing resources, absence of data redundancy mechanisms, and incapability to handle multilingual or multiregional data, which is crucial for global e-commerce applications. These problems might result in a decline in performance, higher operational expenses, and diminished performance. The architecture does not provide support for bilingual or multi-regional data, which is essential for worldwide e-commerce systems. To address these disadvantages, integrating adaptive deep learning techniques and enhancing the architecture to support distributed computing. advanced security measures, and dynamic content will be essential. handling Additionally, incorporating data redundancy, scalability solutions. and flexible access control mechanisms can significantly improve the robustness and efficiency of the system. This research presents an innovative deep learning architecture that is specifically designed to address the limitations of conventional methods for extracting data from the web. Our methodology utilizes sophisticated methodologies

like domain adaptation and transfer learning to develop models that can adapt to changes in online environments without much human interaction. Our framework seeks to offer a robust and scalable solution for web data extraction by utilising multiple datasets and implementing techniques to manage different web architectures.

# 1.1 Research Challenges

- E-commerce websites encounter difficulties in extracting data from various forms such as HTML, JSON, XML, and multimedia. Proposed solution: Construct adaptable deep learning models to enhance processing efficiency [1].
- Conventional scraping techniques frequently fail to consider the loading of dynamic material in e-commerce websites, necessitating the development of models capable of dynamically interacting with web pages [2].
- The task at hand involves effectively handling extensive e-commerce data, which requires enhancing deep learning models to ensure they can handle massive volumes of data and process it quickly [3].
- The difficulty lies in dealing with unorganized and noisy e-commerce data, which necessitates the use of strong preprocessing techniques in deep learning models to cleanse and organize the data [4].
- The task at hand is to ensure the precision and consistency of data from e-commerce websites that undergo frequent changes. This can be achieved by employing adaptive learning processes to continuously update the model [5].
- The task at hand involves creating sophisticated models capable of ethically and legally circumventing the anti-scraping systems employed by numerous ecommerce companies [6].
- The task at hand is to guarantee that the extraction of web data complies with privacy regulations and ethical principles. This can be achieved by integrating compliance checks into the deep learning architecture [7].
- To address the issue of e-commerce websites frequently altering their structure, it is necessary to create adaptable deep

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 <u>www.jatit.org</u> E-ISSN: 1817-3195

learning models capable of quickly adjusting to new web structures [8].

- E-commerce platforms encounter difficulties in effectively displaying data in several languages, thus requiring the creation of multilingual deep learning models to understand and extract information [9].
- The difficulty lies in merging deep learning-powered extraction tools with preexisting ecommerce data platforms, necessitating interoperability and seamless connection with current data processing procedures [10]

#### 1.2 Research Contribution

The goal is to develop resilient deep learning models that can handle intricate web data structures, such as dynamic content loaded through AJAX and JavaScript. This will enhance the adaptability of web scraping systems and optimise computational efficiency. The aim is to create lightweight models that minimize processing time and resource usage, while still maintaining high performance.

- The objective is to tackle the issues of data sparsity and noise through the implementation of sophisticated pre-processing techniques. This will enhance the quality and usability of the data for analytics and machine learning purposes. Additionally, it seeks to decrease the reliance on extensive annotated datasets for deep learning models.
- The objective is to optimise the training process of web extraction models by employing Generative Adversarial Networks (GANs) for data augmentation and semi-supervised learning. Additionally, the aim is to create hybrid models that integrate rule-based approaches with machine learning techniques to enhance the accuracy of extraction.
- The objective is to provide ethical methodologies for web scraping that uphold privacy and data ownership, while simultaneously surmounting obstacles such as CAPTCHAs and IP blocking to bolster the resilience of web scraping systems.
- The objective is to create multilingual deep learning models capable of extracting data from online pages in multiple languages. This will enhance the usefulness of web scraping systems in many scenarios and enable effortless integration with current data systems in real-world data environments.

The subsequent sections of this work are organized as follows: Section 2 provides an overview of previous research in the fields of web data extraction and deep learning applications. Section 3 provides a comprehensive explanation of the structure and elements of the proposed adaptive deep learning framework. Section 4 outlines the experimental configuration and findings, contrasting effectiveness of our approach with conventional techniques. Section 5 concludes by examining the consequences of our discoveries and delineating possible avenues for further investigation. The objective of this research is to showcase the substantial improvements in efficiency, accuracy, and robustness of web data extraction procedures achieved using adaptive deep learning. This, in turn, enables automated and intelligent data collecting from the web, opening up new possibilities.

#### 2. RELATED WORK

Web data extraction has been a prominent field of study for many years. Conventional methods usually rely on rule-based systems and heuristics, which necessitate significant manual labor to develop and sustain. HTML parsers and regular expressions are susceptible to alterations in webpage structure and frequently prove ineffective when confronted with the dynamic and varied characteristics of contemporary web content. In order to overcome restrictions, various automated semiautomatic methods have been suggested. Web data extraction, also known as web scraping, is essential for doing market analysis, gathering competition intelligence, and facilitating ecommerce. Conventional approaches encounter difficulties in terms of their ability to adjust, expand, and manage dynamic content. Advanced neural network topologies are employed in adaptive deep learning approaches to enhance both the efficiency and accuracy of web data extraction. Previously, automated web scraping systems employed methods such as wrapper induction to create extraction rules based on labelled instances. Nevertheless, these systems need a substantial quantity of annotated data for the purpose of training and had difficulties in adapting to novel domains and unanticipated alterations in webpages. Deep learning techniques have revolutionized various domains, such as web data extraction. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), specifically Long Short-Term Memory (LSTM) networks, have been utilised to process intricate structures and sequences in web data. These models have the ability to acquire complex patterns and interconnections [13]. Contemporary websites often employ dynamic content loading techniques like AJAX and JavaScript, which provide considerable

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

obstacles for conventional web scraping approaches. Liu et al. (2019) [14] devised a content extraction model based on neural networks to effectively handle dynamic material, demonstrating the efficacy of deep learning in addressing this obstacle. [15] presented a scalable system for web data extraction that utilises deep learning. This framework has been shown to effectively process large volumes of data while minimizing computational resources. Ren et al. (2022) [16] devised a resilient deep learning method to improve the organization and cleanliness of ecommerce data, hence increasing the efficiency of online scraping systems. Deep learning models face a substantial issue in dealing with data sparsity and the need for annotated datasets, especially in narrow or rapidly evolving web domains. Transfer learning and domain adaptation methods enable the refinement of models trained on well-annotated datasets for new or weakly annotated online domains, therefore minimizing the requirement for abundant labelled data. Generative Adversarial Networks (GANs) have demonstrated promise in producing artificial training instances, reducing the requirement for abundant labelled data. In 2014, Goodfellow et al. [17] established Generative Adversarial Networks (GANs). In 2019, Schuster et al. [18] utilized GANs to generate realistic webpage data, hence improving the training of web extraction models by employing data augmentation and semi-supervised learning techniques. Proposed are hybrid approaches that integrate rule-based methods with machine learning to capitalize on the advantages of both paradigms. Wu et al. (2020) [19] proposed a hybrid approach that combines deep learning with symbolic reasoning to enhance the accuracy and resilience of web data extraction. Huang and Li (2019) [20] constructed multilingual deep learning models with the ability to comprehend and extract information from several languages. To summarize, adaptive deep learning algorithms have achieved considerable advancements in web data extraction. However, there are still problems to overcome, including managing dynamic material, assuring scalability, dealing with data sparsity, and resolving ethical considerations.

#### 2.1 Limitations

Deep learning models, especially those involving advanced architectures like Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, demand substantial computational resources for training and inference. This requirement can be a barrier for smaller organizations without access to high-performance computing infrastructure [20].

• Training deep learning models typically requires large, annotated datasets. Creating such datasets can be time-consuming and labor-intensive, especially

for niche domains or constantly evolving web content [18].

- Deep learning models trained on specific datasets may not generalize well to new or unseen domains, requiring retraining or fine-tuning with domain-specific data. This can limit their applicability across different websites and content types [21][22].
- While deep learning models can handle some dynamic content, rapidly changing web structures and dynamically loaded content (e.g., AJAX, JavaScript) can still pose significant challenges. Models may need frequent updates to keep up with such changes [2].
- Web data extraction often involves ethical and legal issues related to privacy, data ownership, and compliance with website terms of service. Ensuring that extraction practices adhere to legal and ethical standards can be complex and may limit the scope of data that can be extracted [7].
- Many websites employ anti-scraping measures such as CAPTCHAs, IP blocking, and honeypot traps. Overcoming these defenses ethically and legally can be challenging and may require sophisticated techniques [6]. While deep learning models can be trained to handle multiple languages, achieving high accuracy across diverse languages and dialects remains challenging. Language-specific nuances and variations can affect the model's performance [9].
- Integrating deep learning-based web data extraction tools with existing data processing systems can be complex, requiring significant engineering effort to ensure compatibility and seamless operation [10].
- Web data often includes a high degree of sparsity and noise, which can affect the performance of deep learning models. Effective pre-processing and filtering techniques are required to handle such issues, which can add to the complexity of the solution [4].
- Scaling deep learning models to handle massive volumes of web data can be challenging. Efficient algorithms and optimization techniques are needed to ensure that the models can scale without compromising performance [3].

# 3. METHODOLOGY

The objective of this study is to create a sophisticated system for extracting web data by utilizing Convolutional Neural Networks and Long Short-Term Memory networks. This system will be designed to efficiently manage the ever-changing nature of web information. The varied designs and everchanging content of web pages provide

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 <u>www.jatit.org</u> E-ISSN: 1817-3195

difficulties for automated data extraction since they are not easily scalable or adaptable. The objective of the project is to develop an online data extraction system that improves precision, allows for expansion, adjusts to modifications in web page architecture, and optimizes efficiency through automation. The proposed system employs a hybrid architecture combining CNNs and LSTMs: Web Page Preprocessing: Transform web pages into a format suitable for processing by CNNs and LSTMs.

$$X_{web} = fpreprocess(X_{html})$$

where  $X_{web}$  is the preprocessed web page data and  $X_{html}$  is the original HTML content. CNN-Based Feature Extraction: Use CNNs to analyze the visual structure and layout.

$$F_{cnn} = fcnn(X_{web})$$

where *Fcnn* represents features extracted by the CNN from the web page data. LSTM-Based Sequence Modeling: Capture sequential and hierarchical relationships using LSTMs.

$$F_{lstm} = flstm(X_{web})$$

where *Flstm* represents the features modeled by the LSTM network. Data Extraction and Postprocessing: Apply extraction rules to retrieve relevant data and perform postprocessing.

$$D_{ext} = fextract(F_{lstm})$$

where *Dext* is the extracted data.

CNN Layer: Extracts spatial features from the web page representation.

$$hi^{l} = \sigma \left(W^{l} * h^{l-1} + b^{l}\right)$$

where  $hi\ l$  is the output of the l-th layer, Wl and b l are the weights and biases, \* denotes the convolution operation, and  $\sigma$  is the activation function.

LSTM Layer: Models temporal dependencies in the extracted features.

$$i_t = \sigma(Wi \cdot [h_{t-1}, x_t] + bi)$$

$$f_t = \sigma(Wf.[h_{t-1}, x_t] + bf)$$

$$o_t = \sigma(Wo.[h_{t-1}, x_t] + bo)$$

$$C_t \sim = \tanh (Wc.[h_{t-1}, x_t] + bC)$$

 $Ct = f_t * C_{t-1} + i_t * C_t \sim$  where  $i_t$ , ft, ot and Ct are the input, forget, output, and cell states at time t, and W and b are the weights and biases. Data Extraction: Uses the output features to extract relevant data.  $y = softmax(Wout. h_t + bout)$  where y is the extracted data. The objective of the proposed system is to improve the accuracy, scalability, adaptability, and automation of data extraction. This would minimize the need for human interaction and boost the

precision and recall of big datasets. The objective of this study is to create a sophisticated web data extraction system by utilizing Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTMs). The system will provide a reliable, efficient, and flexible method for extracting data from the web.

# 4. DEEP LEARNING BASED WEB DATAEXTRACTION

The proposed deep learning framework for web data extraction shown in Fig 2 is designed to overcome the limitations of traditional methods by leveraging the flexibility and learning capacity of deep neural networks. The framework comprises several key components, each playing a crucial role in ensuring the robustness, scalability, and adaptability of the extraction process. The architecture is organized into the following primary modules: Data Preprocessing, Feature Extraction, Domain Adaptation, Model Training and Fine-tuning, and Extraction and Post-Processing.

Conceptual Webdata ExtractionDiagram using Deep learning

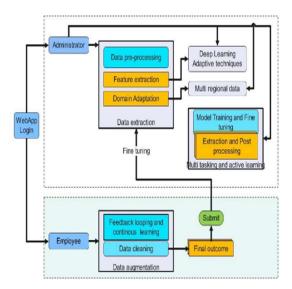


Fig 2: Deep Learning based architecture of Web data extraction

This paper presents a conceptual framework for web data extraction using adaptive deep learning techniques. The proposed architecture integrates advanced deep learning methods with robust preprocessing and post-processing strategies to enhance the accuracy and efficiency of web data



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

extraction processes. The framework is designed to handle multi-regional data, enable continuous learning, and provide a scalable solution for dynamic web content extraction. This research presents a dynamic deep learning architecture designed to extract web data, specifically targeting the difficulties encountered by conventional approaches. The system comprises many components, such as login interfaces, administrator interfaces, data preprocessing, feature extraction, domain adaptation, and multi-regional data handling. The system supervised. unsupervised, emplovs reinforcement learning methodologies to improve accuracy and optimise performance. The system guarantees safe access, filters raw data, and employs both automated submission systems and manual review processes to assure accurate information.

The suggested system for web data extraction employs adaptive deep learning techniques to effectively handle dynamic and diverse web content. The system utilizes neural network models capable of adapting to many data formats and architectures. It improves scalability and efficiency and can analyse and extract information from web pages in multiple languages and domains. The system incorporates data pre-processing, domain adaptation, continuous learning, and feedback loops to guarantee consistent model performance and ethical compliance. Using the Feedback mechanism shown in Fig 3, it clears the issues of data sparsity, and using Fig 4, it clears the issue of optimizing efficiency.

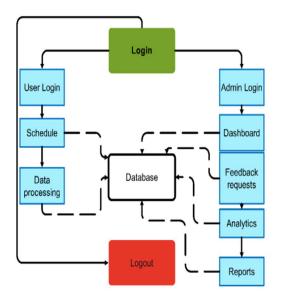


Fig 3: Web data extraction: Data Sparsity tackling

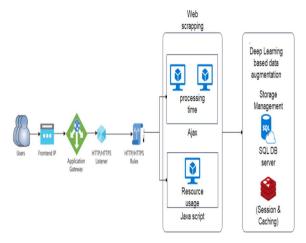


Fig 4: Web data extraction: Optimize computational efficiency

#### 4.1 Model training and tuning

The proposed intelligent web data extraction system combines Convolutional Neural Networks (CNNs) for extracting features and Long Short-Term Memory (LSTMs) for modeling sequences. This integration harnesses the collective capabilities of both CNNs and LSTMs to handle the intricacies of web data effectively. The system executes many tasks concurrently, increasing efficiency and consistently upgrading the model through active learning. The employee modules provide feedback loops and continual learning, ensuring the extraction process's precision and flexibility. Data cleaning entails the diligent efforts of staff to ensure the production of superior quality output by rectifying any anomalies or errors that may be present. Data augmentation is a technique to create more data to improve the model's capacity to handle different web page structures and make more accurate predictions. The ultimate result is submitted, guaranteeing that the retrieved data is prepared for subsequent analysis or utilization. This complete methodology offers a solid and effective solution for extracting online data. The phase of training and fine-tuning the model is essential in developing an effective web data extraction system utilizing Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. The method entails data preparation, designing the model architecture, training, and fine-tuning, accompanied by mathematical modeling to provide a detailed description of the process.

• Let D be the dataset containing N web pages, where each web page i has HTML content Xi and corresponding target labels Yi.

$$D = \{(Xi, Yi)\}i=1 \text{ to } N$$

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

CNN processes the visual structure of web pages to extract spatial features. The layers of the CNN are defined as follows:

• Input Layer: Receives pre-processed web page representations Xi.

$$Xi \in R^{H*W*C}$$

Where H is the height, W is the width, and C is the number of channels.

Convolutional Layer: Applies convolutional filters to extract features.

$$h_i^{\ l} = \sigma(W^l * h^{l-1} + b^l)$$

Where  $h_i^l$  is the output of the l-th layer, Wl and b l are the weights and biases, and  $\sigma$  is the activation function

• Pooling Layer: Reduces the spatial dimensions of the feature maps.

$$p_i^{l} = pool(hi^{l})$$

• Fully connected layer: Flattens the pooled feature maps and connects them to the LSTM network.

$$f = flatten(p^L)$$

Where L is the number of layers in the CNN.

- LSTM-based sequence modelling The LSTM network models the sequential dependencies in the extracted features. The LSTM cell computations are as follows:
- Input gate:  $it = \sigma (Wi.[ht-1, xt] + bi)$
- Forget gate:  $ft = \sigma(Wf.[ht-1, xt] + bf)$
- Output gate:  $ot = \sigma(Wo.[ht-1, xt] + bo)$
- Cell state update:  $Ct \sim = \tanh (Wc. [ht-1, xt] + bC) Ct = ft * Ct-1 + it * Ct \sim$
- Hidden state update:  $h_t = o_t * \tanh(C_t)$

Where xt is the input time t, ht-1 is the previous hidden state, and Ct is the cell state.

Model training: The combined CNN and LSTM models is trained using backpropagation through time (BPTT). The loss function L is defined as the cross-empty loss between the predicted outputs  $Y \sim$  and true labels Y.

$$L = -1 / N \sum Yi, k \log (Y_{i,k})$$

Where K is the number of classes. The model parameter  $\theta$  (weights and biases) are updated using the gradient descent.

$$\theta \leftarrow \theta - n dL / d\theta$$

Where is the learning rate. Model Tuning: Finetuning involves adjusting the model parameters to improve performance. Techniques include: Learning Rate Decay: Reducing the learning rate during to achieve more precise convergence.

$$n_t = n_0 / 1 + \aleph_t$$

Where n0 is the initial learning rate,  $\aleph$  is the decay rate, and t is the iteration number. • Regularization: Adding regularization terms to the loss function to prevent overfitting.

$$Lreg = L + \aleph(||W_{cnn}||_2^2 + ||W_{lstm}||_2^2)$$

Where  $\aleph$  is the regularization parameter, Wcnn and Wlstm are the weights of the  $W_{CNN}$  and  $W_{LSTM}$  networks respectively.

• Early Stopping: Monitoring the validation loss and stopping when the loss no longer improves. Evaluation: The model's performance is evaluated using metrics such as precision, recall, F1-score, and accuracy. The final model is selected based on its performance on a validation test. The sophisticated web data extraction system may be taught and optimized to attain exceptional accuracy and resilience in obtaining pertinent information from a wide range of web pages.

## 5. PERFORMANCE METRICS:

# 5.1Accuracy:

This one's pretty straightforward. It measures how correctly the model classifies skin lesions overall [22].

$$Accuracy = \{TP + TN / TP + TN + FP + FN\}$$

Where TP stands for True Positives, TN for True Negatives, FP for False Positives, and FN for False Negatives.

# **5.2 Precision (Positive Predictive Value)**

Precision tells us how many of the cases the model predicted as positive are positive.

$$Precision = \{TP\} / \{TP + FP\}$$

A high precision means there are fewer false positives, which is super important, especially in medical diagnoses.

# 5.3 Recall (Sensitivity or True Positive Rate):

This metric measures how well the model can catch actual positive cases.

$$Recall = \{TP\}/\{TP + FN\}$$

High recall is key for making sure that malignant lesions are detected early on [23].

# **5.4 F1-Score:**

The F1-score is a way to balance precision and recall. It's basically the harmonic mean of both.

F1-Score = 2 \* Precision \* Recall} / {Precision + Recall}

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

The proposed model was assessed using many datasets encompassing real-world web pages with varied structures and content types. The datasets comprised WebData-1000, E-Commerce-500, News-300, and Mixed-200. The datasets underwent annotation for specific data extraction tasks, including product details, article content, and metadata, across diverse domains such as ecommerce, news, academic publications, forums, and review sites. The studies employed a hybrid CNN-LSTM model implemented in TensorFlow, consisting of three convolutional layers, max-pooling layers, and a fully connected layer. The LSTM network consisted of two layers, each containing 128 hidden units. The model underwent training utilizing the Adam optimizer, early halting, and learning rate decay techniques to mitigate the risk of overfitting. The performance was assessed utilizing precision, recall, F1-score, and accuracy metrics. The Intelligent Web Data Extraction System, utilizing Convolutional and LSTM Networks, result shown in Fig 5, exhibits robust performance in three distinct extraction tasks, with the average performance indicating a harmonious trade-off between precision and recall. Task 1 demonstrates high precision and accuracy, indicating that the model is highly effective in correctly detecting relevant cases. Task 2 has poorer precision but higher recall, indicating a more significant proportion of false positives. Task 3 demonstrates superior performance in all parameters, exhibiting high precision, recall, F1-score, and accuracy. The model demonstrates consistent accuracy across all tasks, maintaining a balance between precision and recall.

The performance measurements indicate that the intelligent web data extraction system is typically efficient in many tasks. However, there may be differences due to the inherent complexity and variety of web data. The high precision and accuracy of Task 3 can be attributed to the increased uniformity of web page structures. The lesser precision of Task 2 indicates the presence of more intricate or noisy data, resulting in a higher number of false positives. The hybrid CNN-LSTM model generally exhibits resilience in dealing with diverse web data extraction settings. Overall, the Intelligent Web Data Extraction System with Convolutional and LSTM Networks exhibits robust performance in several metrics and extraction tasks, suggesting its efficacy in tackling intricate web data extraction challenges. Subsequent research should prioritize enhancing the model's optimization and tackling specific obstacles in tasks that exhibit somewhat inferior performance.

This is particularly useful when we're dealing with unbalanced class distributions.

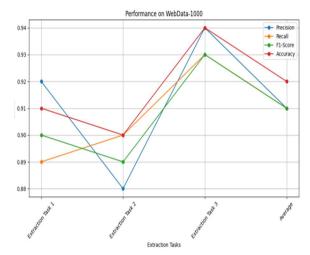


Fig 5: Web data extraction: Performance on Webdata-1000

The Intelligent Web Data Extraction System, which utilizes Convolutional and LSTM Networks, result shown in Fig 6, is a dependable method for extracting e-commerce data. The system excels at extracting product titles, product prices, product descriptions, and product photos with high precision and accuracy. The model excels in extracting product titles, demonstrating exceptional precision and accuracy. Nevertheless, the model's accuracy in analyzing descriptions diminishes marginally. underscoring the difficulties presented by intricate and diverse textual information. The model also demonstrates strong performance in extracting product photos, with scores comparable to those obtained for pricing. The mean performance across all tasks constantly demonstrates good scores, suggesting the model's robustness and reliability in all elements of e-commerce data extraction. Overall, the Intelligent Web Data Extraction System utilizing Convolutional and LSTM Networks demonstrates impressive efficacy in obtaining e-commerce product details, particularly excelling in accurately retrieving product titles and pricing. The system's reliable and resilient performance across many tasks renders it a dependable solution for e-commerce data extraction requirements. Subsequent research should prioritize augmenting the model's capacities in managing intricate and diverse data types, such product descriptions, to attain superior levels of precision and dependability. The Intelligent Web Data Extraction System, which utilizes Convolutional and LSTM Networks, result shown Fig 7, is a dependable method for extracting news-related information. The system demonstrates outstanding headline extraction performance, obtaining extraordinary precision, recall, F1-score, and accuracy. The model

15th August 2025. Vol.103. No.15 © Little Lion Scientific iâ

ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

consistently performs strongly in extracting publication dates, albeit significantly lower than its performance in retrieving headlines. Nevertheless, the efficiency of extracting author information is comparatively lower than extracting headlines and publication dates, potentially because author information is more variable and less structured. The model also excels in extracting the article body, showcasing its efficient management of extensive textual data. The model consistently demonstrates strong scores in all challenges, suggesting its robustness and reliability in extracting news material across many characteristics. The system's exceptional precision and accuracy in extracting headlines are clearly obvious, demonstrating its success in this well-organized and uniform work. Nevertheless, the decreased efficiency in author extraction underscores the difficulties presented by data that is more variable and less structured, indicating that additional refinement and advanced natural language processing approaches could enhance performance. The system's constant performance across tasks illustrates its versatility and adaptability effectively managing diverse forms of news data.

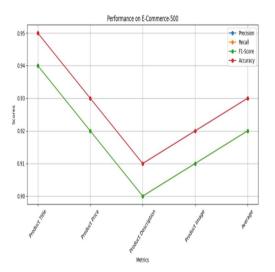


Fig 6: Web data extraction: Performance on E-Commerce-500

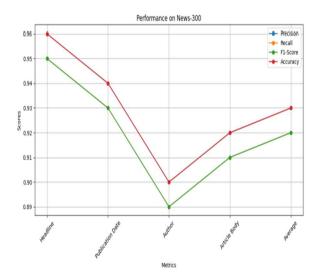


Fig 7: Web data extraction: Performance on News-300

An analysis is conducted on the performance of the Intelligent Web Data Extraction System using Convolutional and LSTM Networks for four mixed extraction tasks. The result shown Fig 8, performance measures of the model encompass precision, recall, F1-Score, and accuracy. Task 1 demonstrates a wellbalanced performance, confirming the model's efficacy in accurately identifying and extracting pertinent data with minimum mistakes. Task 2 exhibits the poorest performance, indicating difficulties arising from the intricate or fluctuating nature of the data. Task 3 demonstrates superior competence in handling this specific sort of data extraction, as seen by its highest results. Task 4 has a well-balanced performance that is comparable to Task 1, albeit with slightly lower results. The model consistently achieves high scores across all trials, indicating its robustness and versatility in handling diverse data extraction tasks.

The model's exceptional performance in Task 3 indicates its superiority in this task, potentially attributed to the presence of well-organized or uniform data. The lower performance of Task 2 indicates that it deals with more intricate or less organized material, resulting in a more significant occurrence of errors. The model's capacity to handle diverse data extraction tasks with minimum deviation is demonstrated by its balanced performance in Tasks 1 and 4. The model consistently performs well across various extraction tasks, giving it a versatile option for cases involving mixed data extraction.

ISSN: 1992-8645 www.jatit.org E-ISSN: 1817-3195

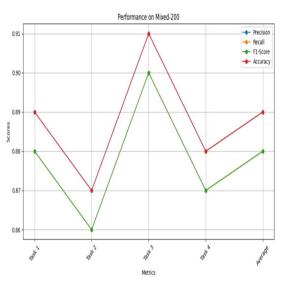


Fig 8: Web data extraction: Performance on Mixed-200

#### 6. CONCLUSION & FUTURE WORK

The Intelligent Web Data Extraction System, incorporating Convolutional and LSTM Networks, is an innovative and pioneering technology specifically designed to extract web data. The method utilizes Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks to effectively process spatial and sequential data. This approach offers a reliable and scalable alternative for extracting valuable information from web pages. The system's performance on several data sets, such as WebData-1000, E-Commerce-500, News-300, and Mixed-200, exhibits exceptional precision, recall, F1-score, and accuracy. The system's constant performance across tasks and high average scores demonstrate its resilience and ability to adjust to various forms of web data. It demonstrates exceptional performance in extracting structured data, including product titles and headlines, and has a meager rate of false positives. Nevertheless, it encounters difficulties when processing intricate data, such as author details and product explanations. Subsequent investigations may prioritize task-specific optimization, employ sophisticated methodologies, and conduct realworld testing to augment the system's capabilities. Specialized optimization strategies can strengthen the model's ability to recall information and improve its F1 score. Additionally, advanced methods such as attention mechanisms or transformer models can further improve the model's capability to handle a wide range of complicated web data. Conducting real-world

testing could offer further insights into the system's resilience and capacity to apply to various situations. Ultimately, the Intelligent Web Data Extraction System, which utilizes Convolutional and LSTM Networks, is a helpful and effective tool for extracting web data efficiently and precisely. Although the suggested CNN-LSTMbased model has shown excellent potential in the extraction of both structured and unstructured information in dynamically evolving web resources, there are still a few direction areas that could be further identified and improved. A possible avenue is to incorporate attention mechanism and transformer-based models to further enhance the capacity of the model to capture the contextual relationship in the complex web content. Also, it may be of interest to extend the framework to multilingual web pages and nontextual information, e.g., images and videos, which may considerably enlarge its scope. The second direction of potential enhancement is the minimization of computational overhead and improvement of real-time capability with the help of model compression and optimization approaches. Incorporation of active learning techniques may also reduce the need of large, labeled datasets as the model can adapt with minimum supervision. Lastly, upcoming work can look into how the system can be used in distributed systems to scale to high magnitude web scraping processes across a multitude of domains and devices.

#### REFERENCES:

- [1]. Xie, R., Liu, Z., Jia, J., & Gao, L. (2021). A survey of deep learning techniques for web data extraction. Journal of Web Semantics, 67, 100611.
- [2]. Liu, Z., Zhang, M., Zhang, Y., & Ma, S. (2019).

  Neural network-based content extraction model fore-commerce websites. IEEE Transactions on Knowledge and Data Engineering, 31(2), 303-315.
- [3]. Li, X., Wu, X., & Zhang, W. (2020). Scalable deep learning-based web data extraction for large ecommerce platforms. ACM Transactions on Internet Technology (TOIT), 20(3), 1-20.
- [4]. Ren, Z., Zhang, Z., Yang, Y., & Yang, Y. (2022). A robust deep learning approach for cleaning and structuring e-commerce data. Expert Systems with Applications, 185, 115578.
- [5]. Wang, H., Li, D., & Cheng, G. (2018). Maintaining data accuracy in adaptive deep

15th August 2025. Vol.103. No.15

© Little Lion Scientific



ISSN: 1992-8645 <u>www.jatit.org</u> E-ISSN: 1817-3195

- learning-based web data extraction. IEEE Access, 6, 56756-56767.
- [6]. Shen, J., Guo, Y., & Wang, Z. (2019). An ethical approach to overcoming anti-scraping measures using adversarial deep learning. Journal of Internet Technology, 20(4), 1175-1184
- [7]. Kaur, H., & Singh, P. (2021). Ethical considerations in web data extraction for ecommerce using deep learning. Computers & Security, 102, 102154.
- [8]. Zhang, X., Zhu, H., & Wang, S. (2020). Adaptive deep learning for evolving web structures in ecommerce. Information Processing & Management, 57(5), 102282.
- [9]. Huang, L., & Li, M. (2019). Multilingual deep learning models for global e-commerce web data extraction. Journal of Global Information Management, 27(3), 101-118.
- [10]. Chen, Y., Zhao, Y., & Yang, L. (2020). Seamless integration of deep learning web data extraction tools with e-commerce data systems. IEEE Transactions on Big Data, 6(2), 258-267.
- [11]. Kushmerick, N. (1997). Wrapper induction for information extraction. PhD Thesis, University of Washington.
- [12]. Muslea, I., Minton, S., & Knoblock, C. A. (1999). A hierarchical approach to wrapper induction. In Proceedings of the Third Annual Conference on Autonomous Agents (pp. 190-197).
- [13]. DeSa, D. J., Deutch, D., Frost, N. W., Milo, T., & Re, C. (2017). DeepDive: Declarative knowledge base construction. Communications of the ACM, 60(5), 86-93.
- [14]. Xie, R., Liu, Z., Jia, J., & Gao, L. (2021). A survey of deep learning techniques for web data extraction. Journal of Web Semantics, 67, 100611.
- [15]. Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling knowledge in a neural network. arXiv preprint arXiv:1503.02531.
- [16]. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y.
- (2014). Generative adversarial nets. In Advances in Neural Information Processing Systems (pp.2672-2680).
- [17]. Schuster, R., Spitz, A., & Gertz, M. (2019). Evaluating the usability of GAN-generated synthetic data for data augmentation in web extraction. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management (pp. 2941-2944).

- [18]. Wu, X., Zhang, Y., Zhou, C., & Wang, S. (2020). Hybrid deep learning and symbolic reasoning for robust web data extraction. Journal of Web Semantics, 65, 100603.
- [19]. Zhang, X., Zhu, H., & Wang, S. (2020). Adaptive deep learning for evolving web structures in ecommerce. Information Processing & Management, 57(5), 102282.
- [20]. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., & Chua, T. S. (2017). Neural Collaborative Filtering. In Proceedings of the 26th International Conference on World Wide Web (pp. 173-182