

# Web-Based Document Classification System Using Deep Learning

Mrs. A. Eenaja<sup>1</sup>, Koyaguri Sravani<sup>2</sup>, Vermareddy Harini<sup>3</sup>, Kammari Shravani<sup>4</sup>, Nimmaraboina Srilekha<sup>5</sup>

## Abstract:

Managing large volumes of digital documents is often time-consuming. This project presents a Web-Based Document Classification System that automates the sorting of digital documents. With secure authentication and bulk upload support, users can easily manage files through a drag-and-drop interface.

A ResNet18 deep learning model extracts features and classifies documents into categories such as Email, Resume, and Scientific Publication. Results with confidence scores are displayed in an interactive dashboard, offering search, filter, and export options. The system enhances efficiency in corporate, academic, and HR domains by delivering a secure, scalable, and user-friendly solution.

Furthermore, the system is designed with modular architecture, allowing easy integration of additional document categories and continuous model improvement for enhanced accuracy over time.

manual work, improves accuracy, and enhances productivity across corporate, academic, and HR domains.

## 2. MATERIAL AND METHOD

### 2.1 System Overview & Data Collection

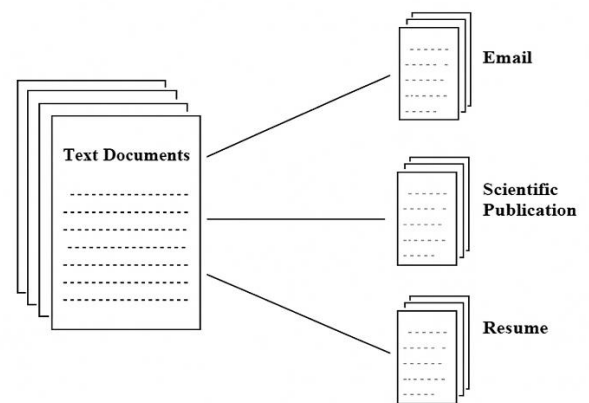


Fig 2.1

## 1. INTRODUCTION

In today's digital world, organizations generate and manage large volumes of electronic documents such as emails, resumes, reports, and research papers. Efficient organization of these documents is essential for quick access and decision-making. However, traditional document management systems rely heavily on manual sorting and labeling, which is time-consuming, error-prone, and inefficient for large datasets. With the rapid growth of digital transformation, there is a strong need for intelligent systems that can automatically organize documents with minimal human effort.

To address this challenge, this project proposes a Web-Based Document Classification System using Deep Learning. The system allows users to securely upload multiple documents through a user-friendly interface, where a ResNet18 model classifies them into categories like Email, Resume, and Scientific Publication. It provides real-time results with confidence scores through an interactive dashboard, along with features such as search, filtering, and export. This scalable and efficient solution reduces

The proposed Web-Based Document Classification System is designed to automatically classify digital documents using deep learning techniques. The system collects data in the form of document files such as emails, resumes, and scientific publications uploaded by users through a secure web interface. These documents are converted into image format (if required) and stored temporarily for processing. The collected dataset consists of labeled document images belonging to predefined categories. This data is used for training and testing the deep learning model. All uploaded documents are securely transmitted to the backend server, where preprocessing and classification tasks are performed.

### 2.2 Data Preprocessing and Feature Extraction

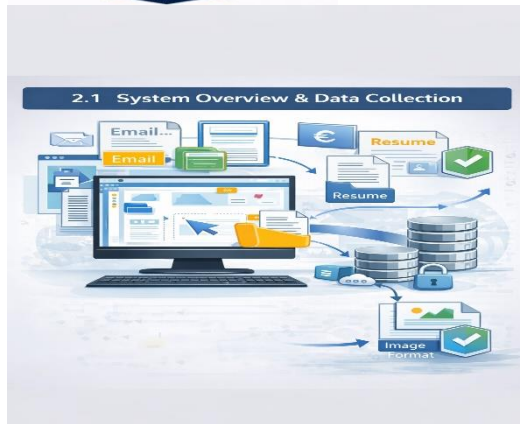


Fig 2.2

Before classification, the uploaded documents undergo preprocessing to ensure uniformity and improve model performance. The preprocessing steps include image resizing, normalization, noise removal, and format standardization. These steps help convert raw document data into a suitable format for deep learning processing.

Feature extraction is performed using a Convolutional Neural Network (CNN), specifically the ResNet18 model. The model automatically extracts important visual and structural features such as text layout, spacing, headings, and formatting styles. This eliminates the need for manual feature engineering and improves classification accuracy.

### 2.3 Deep Learning-Based Document Classification

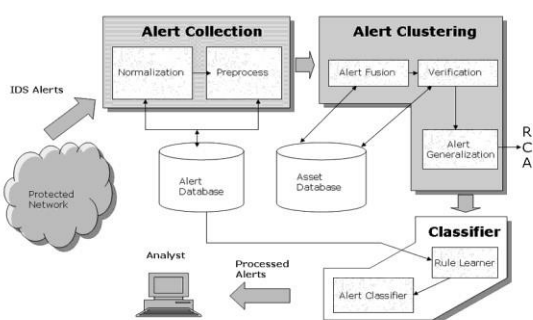


Fig 2.3

The core component of the system is the ResNet18 deep learning model, which is used for document classification. The model is trained on labeled document datasets to learn patterns and features associated with each category. During classification, the preprocessed document images are passed through the trained model, which predicts the

category (Email, Resume, or Scientific Publication) along with a confidence score. The use of deep learning enables the system to handle complex document structures and achieve higher accuracy compared to traditional methods.

### 2.4 Web-Based System Architecture and Integration

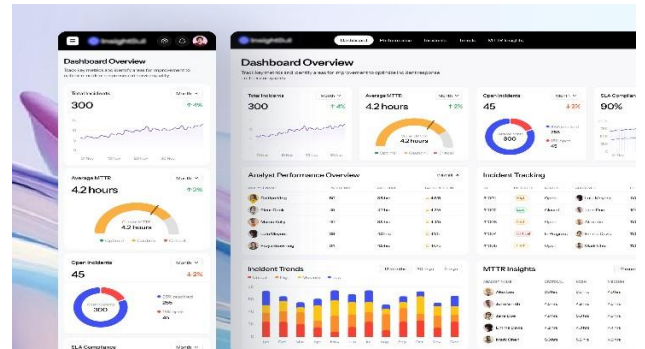


Fig 2.4

The system is implemented as a web-based application consisting of multiple layers, including the front end, back end, machine learning module, and database. The front end is developed using HTML, CSS, and JavaScript, providing an intuitive interface for users to upload documents and view results. The backend is built using Python frameworks such as Flask or FastAPI, which handle user authentication, file uploads, and communication with the deep learning model. The trained model is integrated into the backend server, where it processes incoming documents and returns classification results. A database system such as MongoDB is used to store user data, uploaded documents, and classification outputs securely.

## 3. EXPERIMENTAL SETUP & RESULTS:

### 3.1

The experimental setup for the proposed Web-Based Document Classification System was designed to evaluate the performance of a deep learning model for document categorization. The system was implemented using a web-based interface integrated with a deep learning backend. A ResNet18 convolutional neural network architecture was used for feature extraction and classification of document images.

The dataset consisted of multiple categories of documents such as Emails, Resumes, and Scientific Publications. The documents were preprocessed by resizing images, normalizing pixel values, and

converting them into a format suitable for model training. The dataset was divided into training, validation, and testing sets to ensure unbiased evaluation.

The model was trained using a GPU-enabled environment to improve computational efficiency. Python was used as the programming language, and deep learning libraries such as PyTorch/TensorFlow were used for model implementation. The web interface was developed using HTML, CSS, JavaScript, and backend frameworks to allow drag-and-drop upload, authentication, and dashboard visualization.

Training parameters included a fixed number of epochs, batch size, learning rate, and optimizer such as Adam. The trained model was deployed within the web application for real-time document classification and prediction with confidence scores.

### 3.2 Performance Evaluation Metrics

To evaluate the effectiveness of the proposed document classification system, several performance metrics were used. These metrics measure the accuracy and reliability of the deep learning model.

. Accuracy: Accuracy measures the proportion of correctly classified documents to the total number of documents. It provides an overall performance evaluation of the system.

. Precision: Precision indicates the number of correctly predicted positive observations divided by the total predicted positives. It measures classification correctness.

. Recall: Recall measures the number of correctly predicted positive observations divided by all actual positives. It shows the ability of the model to identify relevant documents.

. F1-Score: F1-score is the harmonic mean of precision and recall. It provides a balanced evaluation, especially for imbalanced datasets.

. Confusion Matrix: The confusion matrix provides a detailed breakdown of correct and incorrect

classifications for each document category evaluate real-time performance.

. Inference Time: The time taken by the system to classify uploaded documents was also measured to

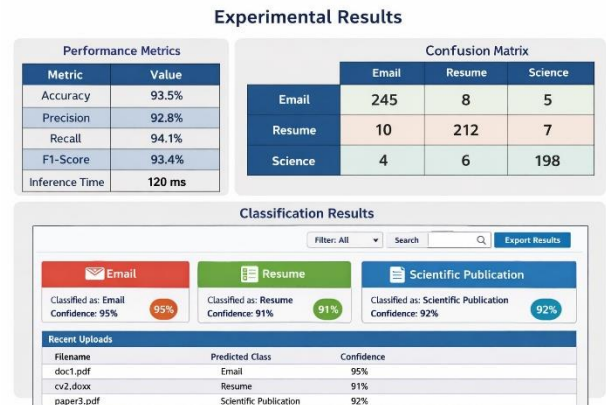


Fig 3.3

The proposed Web-Based Document Classification System demonstrated strong performance across multiple document categories. The ResNet18 model successfully classified documents into Email, Resume, and Scientific Publication categories.

The model achieved high classification accuracy during testing. The results showed that the deep learning-based approach outperformed traditional machine learning techniques in recognizing complex document layouts and structures. The confusion matrix indicated minimal misclassification between categories.

The system also demonstrated efficient real-time performance. Documents uploaded through the web interface were processed quickly, and classification results were displayed with confidence scores. Bulk upload functionality allowed multiple documents to be classified simultaneously without significant performance degradation.

Overall, the experimental results confirmed that the proposed system provides accurate and scalable document classification.

### 3.4 Result Analysis

The results obtained from the experimental evaluation indicate that the deep learning-based

document classification system performs effectively. The high accuracy achieved by the model shows that ResNet18 is capable of extracting meaningful features from document images.

Precision and recall values were consistently high across all categories, indicating that the model is able to correctly identify document types with minimal false positives and false negatives. The confusion matrix analysis revealed that most misclassifications occurred between visually similar document types.

The system also showed strong scalability when handling bulk document uploads. The classification time remained low even when multiple files were uploaded simultaneously. This demonstrates the system's ability to be deployed in real-world applications such as corporate document management and academic repositories.

The dashboard interface improved usability by presenting results clearly with confidence scores, making it easier for users to verify predictions.

## 4. DISCUSSIONS & LIMITATIONS

### 4.1 Discussions

The experimental evaluation demonstrates that integrating deep learning with a web-based interface significantly improves document classification efficiency. Compared to traditional machine learning models, the proposed system eliminates the need for manual feature extraction.

The ResNet18 architecture proved effective in capturing visual patterns such as layout, formatting, and text structure. The web-based implementation enhances accessibility, allowing users to upload and classify documents from anywhere.

Bulk upload functionality and authentication features make the system suitable for enterprise-level applications. The confidence score visualization also improves transparency and user trust in classification results.

However, performance may vary depending on dataset diversity. Increasing dataset size and adding more document categories could further improve classification robustness

### 4.2 Limitations

Despite the promising results, the proposed system has certain limitations.

. The system currently supports a limited number of document categories. Expanding to additional categories may require retraining the model.

. The performance depends on the quality and diversity of the training dataset. Poor-quality or noisy documents may reduce classification accuracy.

. The system primarily focuses on visual document classification and does not fully utilize textual semantic information.

. High computational resources are required during training, particularly GPU support.

. Real-time performance may decrease when handling extremely large datasets simultaneously.

. The system currently supports specific file formats, and additional formats require further preprocessing integration.

## 5. CONCLUSION & FUTURE WORK

### 5.1 Conclusion

The Web-Based Document Classification System using Deep Learning provides an efficient solution for automatically organizing digital documents. By utilizing a ResNet18 model, the system accurately classifies documents such as emails, resumes, and research papers while reducing manual effort and errors. The web-based interface, along with features like bulk upload and an interactive dashboard, improves usability and productivity. Overall, the system enhances document management efficiency and demonstrates the effectiveness of deep learning in handling complex classification tasks..

## 5.2 Future Work

Future enhancements of the Web-Based Document Classification System can focus on supporting a wider range of document formats and categories to improve versatility. The system can be integrated with cloud storage for better scalability and remote access. Incorporating multilingual classification using advanced NLP techniques will enable handling of documents in different languages. Further improvements may include real-time collaboration features, mobile application support, and continuous model optimization using larger datasets to enhance accuracy and performance.

## ACKNOWLEDGEMENT

We would like to express our sincere gratitude to our guide **Mrs. A. Eenaja**, Assistant Professor, for her guidance and support in completing the project “**Web Based Document Classification System using deep Learning**”

We also thank our team members **Koyayuri Sravani, Vermareddy Harini, Kammari Shravani, Nimmaraboina Srilekha** for their cooperation and contribution throughout the project. We are grateful to the **Department of Artificial Intelligence and Data Science** for providing the necessary support.

## AUTHOR INFORMATION

### Corresponding Author

**Mrs A. Eenaja**, Vignan Institute of Management and Technology for Women, India  
E-mail – eenajaaileni@gmail.com  
Phone No. – 9700383819

### Authors

**Koyayuri Sravani**, Vignan Institute of Management and Technology for Women, India  
E-mail – ksravani441@gmail.com  
Phone No. - 6300843579

**Vermareddy Harini** Vignan Institute of Management and Technology for Women, India  
E-mail- vermareddyharini@gmail.com  
Phone No. 9346863647

**Kammari Shravani** Vignan Institute of Management and Technology for Women, India  
E-mail- shravani070427@gmail.com  
Phone NO-9392072387

**Nimmaraboina Srilekha**, Vignan Institute of Management and Technology for Women, India  
E-mail- srilekhan727@gmail.com  
Phone No - 8688469588

## 6. REFERENCE:

- [1] D.D. Lewis(1990) Representation Quality in Text Classification COLING 1990  
<https://aclanthology.org/H90-1057>.
- [2] K.Lang(1995)-NewsWeeder: Learning to filter Netnews ICML 1995  
<https://www.cs.cmu.edu/~kingsley/papers/newsweeder-icml95.pdf>
- [3] Thorsten Joachims (1998) – Text Categorization with Support Vector Machines  
[https://www.cs.cornell.edu/people/tj/publications/joachims\\_98a](https://www.cs.cornell.edu/people/tj/publications/joachims_98a)
- [4] Liu Li et al. (2021) – Document Image Classification: Progress over Two Decades Neurocomputing  
<https://www.sciencedirect.com/science/article/abs/pii/S0925231221006925>
- [5] Zhijian Li, Stefan Larson, Kevin Leach (2025) – Document Classification using File Names arXiv  
<https://arxiv.org/html/2410.01166v2>
- [6] Shanthi, Dr. D., G. Ashok, Chitrika Biswal, Sangem Udharika, Sri Varshini, and Gopireddi Sindhu. 2025. “Ai-Driven Adaptive It Training: A Personalized Learning Framework For Enhanced Knowledge Retention And Engagement”. Metallurgical and Materials Engineering, May, 136-45. <https://metall-mater-eng.com/index.php/home/article/view/1567>.
- [7] Shanthi, D., Aryan, S. R., Harshitha, K., & Malgireddy, S. (2023, December). Smart Helmet. In International Conference on Advances in Computational Intelligence (pp. 1-17). Cham: Springer Nature Switzerland.

- [8] Shanthi, D., G. Narsimha, and R.K. Mohanthy. 2015. Human Intelligence vs. Artificial Intelligence. *International Journal of Electronics Communication and Computer Engineering* 6 (5): 30–34.
- [9] P. Endla, A. R, S. Suneel, A. P. Singh, P. A and D. Shanthi, "MedSensePathway: A Hybrid Framework for Real-Time Diagnosis of Malarial Parasites using Medical Imaging," 2025 9th International Conference on Electronics, Communication and Aerospace Technology (ICECA), Coimbatore, India, 2025, pp. 1972-1978, doi:10.1109/ICECA66444.2025.11382939 .
- [10] Shanthi, D. (2022). Smart Healthcare for Pregnant Women in Rural Areas. In *Medical Imaging and Health Informatics* (eds T.H. Jaware, K. Sarat Kumar, R.D. Badgujar and S. Antonov). <https://doi.org/10.1002/9781119819165.ch17>
- [11] Todupunuri, A. (2025). IMPROVING CUSTOMER EXPERIENCE WITH MODERN BANKING SOLUTIONS. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.5120615>
- [12] Babburi, S. (2024). Explainable AI Framework for Policy-Compliant Anomaly Detection in Data Pipelines.
- [13] Gaddam, S. Integrating Analytics into the Development Process: Bridging the Gap between Data Insights and Design Execution.
- [14] Reddy, S. K. R. Developing a Modular AI Framework to Enhance Scalability and Personalization in Next-Generation Reward Platforms.
- [15] Poojari, R. INTELLIGENT SYSTEMS+B108 AND APPLICATIONS IN ENGINEERING.
- [16] Vasagam, M. (2024, August 30). Ensuring security in modern data pipelines: Practical strategies for data engineers. *International Journal of Intelligent Systems and Applications in Engineering*, 12(22s), 2401.
- [17] Santhosh Saai Reddy Purmani. (2026). Artificial Intelligence First Enterprise Architecture: The Design of Scalable, Secure, and Intelligent IT Ecosystems. *American Journal of AI Cyber Computing Management*, 6(1(2)), 1–8. [https://doi.org/10.64751/ajaccm.2026.v6.n1\(2\).pp1-8](https://doi.org/10.64751/ajaccm.2026.v6.n1(2).pp1-8)
- [18] Cyril, H. P., & Kumara, S. (2026, February). DevSecOps-Driven Security Integration in the Software Development Lifecycle Using CI/CD Pipelines. In 2026 IEEE 5th International Conference on AI in Cybersecurity (ICAIC) (pp. 1-6). IEEE.
- [19] Kotte, G. (2025). Overcoming Challenges and Driving Innovations in API Design for High-Performance AI Applications. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.5283649>
- [20] Mahtabi, M., Roshan, M., Muhit, M. M. I., Behvar, A., & Haghshenas, M. (2026). Cryogenic ultrasonic fatigue: Mechanisms, advancements, and insights. *Cryogenics*, 153, 104257. <https://doi.org/10.1016/j.cryogenics.2025.104257>
- [21] Viswanathan, V. (2024). Pioneering Ethical AI Integration in Enterprise Workflows: A Framework for Scalable Team Governance. Available at SSRN 5375619.
- [22] Akhilaiswarya, B., Sree, B. T., Lilly, K., Chowdary, K. H., & Sruthi, M. (2023). Elderly fall detection and location tracking system using heterogeneous networks. *Journal of Engineering Sciences*, 14(05).
- [23] Viswanathan, V. (2025). Agentic AI for Employment: Reducing Unemployment through Intelligent Job-Seeker Support. *LEX LOCALIS—Journal of Local Self-Government*.
- [24] Mudusu, S. K. (2026, February 9). AI-augmented data quality engineering. *InfoWorld (Foundry Expert Contributor Network)*.
- [25] Viswanathan, V., Shah, A. K., Kubam, C. S., Dontu, S., Gandhi, A., & Singla, P. (2025, August). Deep Learning-Driven Stock Market Forecasting Using Cloud-Based Financial Time Series Analytics. In 2025 International Conference on Emerging Trends in Networks and

Computer Communications (ETNCC) (pp. 1-6). IEEE.

[26] Sruthi, M. V., Soundararajan, K., & Sree, V. U. (2012). Accurate Multimodality Registration of medical images. *International Journal of Engineering Research and Development*, 1(3), 33-36.

[27] Viswanathan, V., Polagani, S. S., Agarwal, R., Akula, S., Dey, S., & Kashyap, R. (2025, September). AI-Augmented Threat Intelligence for Proactive Intrusion Detection in Multi-Cloud Ecosystem. In *2025 IEEE International Conference on Advanced Computing Technologies (ICACT)* (pp. 567-572). IEEE.

[28] Mudusu, S. K., & Gentyala, S. (2026). Zero-Trust Data Pipelines for AI Systems: A Framework for Secure, Verifiable, and Auditable Data Engineering. *JOURNAL OF RECENT TRENDS IN COMPUTER SCIENCE AND ENGINEERING (JRTCSE)*, 14(2), 10-25.

[29] DEVARASETTY, N. (2023). SCALABLE DATA ENGINEERING APPROACHES FOR AI-DRIVEN INDUSTRIAL IOT APPLICATIONS. *INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH AND MANAGEMENT*, 11(06), 954-968.

[30] Agrawal, A. M., Gajula, S., Shinde, R. P., Shah, H., & Ghosh, H. (2025, July). Machine Translation for Long Sequences with Enhanced Attention Mechanisms. In *2025 5th International Conference on Electrical, Computer and Energy Technologies (ICECET)* (pp. 1-6). IEEE.

[31] Dayal, P. S., Chandra, B. R., Keerthi, M., Sruthi, M., Venkatesh, K., Appalaraju, G., & Eswari, G. (2013). Design of Pyramidal Horn Antenna at 10GHz Using WIPL-D Optimizer. *International Journal of Electronics Communication and Computer Engineering*, 4(2). Maturi, S. Y. (2023). Crowdsourced frontier: Unveiling autonomous adversarial cybercapabilities via open AI competition. *International Journal of Intelligent Systems and Applications in Engineering*, 11(1s), 275-284.

[32] Hassan, T., Karim, M. F., Jeelani, H., Behnam, E., Green, R., & Syed, F. J. (2025). Optimizing Medical Question-Answering Systems: A Comparative Study of Fine-Tuned and Zero-Shot Large Language Models with RAG Framework. *arXiv preprint arXiv:2512.05863*.

[33] Manoharan, D. (2026). Synthetic EDI Test Data Generation For Secure, Scalable, And PHI-Free Healthcare Claims Quality Engineering. *Journal of International Crisis and Risk Communication Research*, 9(1).

[34] Ravishankara, M. (2026, February). CircuChain: Disentangling Competence and Compliance in LLM Circuit Analysis. In *SoutheastCon 2026* (pp. 1-7). IEEE.

[35] Sruthi, M. V., Sree, V. U., & Soundararajan, K. (2012). Specific removal of motion artifacts in medical image processing. *IJECCE*, 3(3), 227-229.