



Meta Optimizer Algorithms for Automated Fine-Tuning of Massive Multimodal Models

Hanan Muhajab^{1*}, Dr.R. Udayakumar², Renuka N³, Ugiloy Yunusova⁴, Sharifjon Mirzoyev⁵, Saidavzal Boboyev⁶

¹Department of Computer Science, College of Engineering and Computer Sciences, Jazan University, Jazan, Saudi Arabia; Engineering and Technology Research Center, Jazan University, Jazan, Saudi Arabia. Email: hmuhab@jazanu.edu.sa, <https://orcid.org/0009-0000-3472-3723>

²Professor & Director, Kalinga University, India. Email: rsukumar2007@gmail.com

³Assistant professor, Department of Artificial intelligence and Machine Learning, Kongu Engineering College, Perundurai, India. Email: renuka66ksr@gmail.com

⁴Department of Management and tourism, Tashkent Institute of Irrigation and Agricultural Mechanization Engineers, National Research University, Tashkent, Uzbekistan. Email: ubolkiboyeva@bk.ru, <https://orcid.org/0009-0005-8200-0446>

⁵Associate Professor, Bukhara State Pedagogical Institute Bukhara, Uzbekistan. E-mail: sharifmirzoyev1998@gmail.com, <https://orcid.org/0009-0006-1537-3514>

⁶Researcher, Samarkand State Medical University, Samarkand, Uzbekistan. E-mail: saidavzalbabaev@gmail.com, <https://orcid.org/0000-0002-3725-128X>

*Corresponding author: Email: hmuhab@jazanu.edu.sa

Abstract

The recent breakthroughs in artificial intelligence have helped to develop massive multi-modal models that can process heterogeneous data like text, images, audio, and video. Despite their impressive results in areas like healthcare diagnostics, autonomous systems, robotics, and vision-language learning, fine-tuning these models is computationally expensive and requires considerable hyperparameter tuning efforts. In this paper, a new meta optimizer framework is proposed for automatic fine-tuning of massive multi-modal models. The main objective is to increase the efficiency of the optimization process, ensure stability, and make the model more scalable. The new framework incorporates advanced features like adaptive learning, dynamic hyperparameter tuning, feedback-based optimization, and parameter-efficient fine-tuning to help achieve the scalability of a multi-modal learning environment. The experimental evaluation of the meta optimizer has been performed on multi-modal benchmark datasets and has been compared with the state-of-the-art optimization techniques such as SGD, Adam, AdamW, RMSProp, and Lion Optimizer. The results showed that the proposed meta optimizer achieved the best accuracy and F1-score values, i.e., 96.4% and 0.95, respectively, with the lowest number of convergence epochs being 16 as opposed to 42 for SGD and 31 for Adam. In addition, the proposed framework was able to reduce the training time to 5.1 hours, along with minimizing the GPU memory requirement to 19GB. It can be seen that the training process is quite efficient computationally and scalable as well.

Keywords: Multi-modal Models, Meta Optimizer, Automated Fine-Tuning, Adaptive Learning, Hyperparameter Optimization, Computational Scalability, Deep Learning Optimization

This is an open access article under CC BY 4.0, allowing unrestricted use with proper attribution, a license link, and indication of any changes made.

1. Introduction

Recent advances in artificial intelligence have led to rapid improvements in the design of massive multi-modal models that can process and integrate various kinds of data such as text, images, audio, and videos [7]. Such models, designed based on transformer and foundation models architecture, have gained remarkable success in applications such as visual question answering, speech recognition, medical diagnosis, autonomous systems, and cross-modal learning [2]. Their capability to produce a unified feature representation of different modalities has enhanced contextual reasoning and decision making in modern AI systems [6][11]. Fine-tuning of such massive

models, however, poses several challenges such as high computational costs, huge parameter sizes, memory requirements, and sensitivity to hyperparameters settings. Traditional optimization methods such as Stochastic Gradient Descent (SGD) and Adam do not scale well when faced with heterogeneous modal interactions and dynamic training environment [9][13]. In addition, manual hyperparameter tuning is computationally inefficient. Meta-optimizer algorithms that can automatically tune hyperparameters in a way to improve learning processes during training have been developed as a promising solution [4][15]. Yet, most existing works address single-modal or small-scale problems, emphasizing the necessity of a computationally efficient meta-optimization approach for automated fine-tuning of massive multi-modal models.

This research work presents a new meta optimizer algorithm that can be used in the automated fine-tuning of large multimodal models. This framework adopts adaptive learning approaches, dynamic hyperparameter optimization, and feedback-based parameter tuning to enhance the efficiency and performance of such models.

- 1) Proposed a meta optimizer capable of exploring high-dimensional spaces of hyperparameters for large multi-modal models effectively.
- 2) Used cross-modal feature sensitivities for improving convergence, generalization, and performance.
- 3) Performed extensive experiments and achieved better accuracy, fast convergence, and reduced computational costs compared to alternative optimization algorithms.

The rest of this paper is structured as follows: Section I: Introduction presents the problem of fine-tuning large multi-modal models. Section II: Related Work gives information about previous studies done in the fields of meta-learning and optimization. Section III: Meta-Optimizer Framework describes the method, which includes an automatic fine-tuning process and adaptive optimization using datasets. Section IV: Results and Discussion evaluates the performance, convergence, and scalability of the method. Section V: Conclusion and Future Work conclude the study and suggests future work in optimizing large multi-modal models.

2. Literature Review

Current research has shed light on the fast development of multimodal models that can deal with text, image, audio, and video data. Suggested an adaptive fine-tuning approach for multimodal models; in addition, suggested an AutoML system based on pre-trained transformers for multimodal optimization [1][3]. The paper showed the benefits of fine-tuned multimodal transformer models in both simulated and actual scenarios [5][16]. Created an automated multimodal learning platform combined with large language models; moreover, stressed the importance of development and scalability issues of large multimodal foundation models [6][8].

Moreover, many scholars have investigated meta-learning and optimization approaches for enhancing multimodal systems. Specifically, the suggested reinforcement learning-based alignment method for large multimodal models was created, whereas a meta-learning-based multimodal system for adaptive fault diagnosis [10][12]. In addition, discussed the contribution of nature-inspired meta-heuristic algorithms for solving difficult optimization problems [14][17]. However, some issues regarding scalability, convergence stability, computational overhead, and automated fine-tuning still need to be considered. Consequently, this research suggests a scalable meta optimizer system for automated fine-tuning of multimodal models.

Though there have been improvements in multimodal learning and optimization, issues related to scalability, convergence, and computational efficiency are yet to be sorted out. As a result, an adaptive meta optimizer architecture will be required for effective automated tuning of multimodal models.

3. Proposed Meta Optimizer Architecture

This section describes the proposed meta optimizer architecture that will facilitate automated fine tuning of multimodal models. This proposed system will be able to enhance training efficiency, convergence stability, adaptive learning performance, and reduce computational overheads and human intervention. The proposed system will combine intelligent optimization algorithms and dynamic parameter tuning techniques to enable scalable multi-modal learning of textual, visual, audio, and video data.

3.1 Architecture of the Proposed System

The suggested framework includes five major components which include multi-modal input processing, feature fusion layer, adaptive meta-optimizer engine, dynamic hyperparameter controller, and performance feedback module. First, heterogeneous inputs from different modalities are processed and transformed into common feature representation. The feature representations are then fused through cross-modal attention transformers to produce contextual embeddings. The adaptive meta-optimizer engine continuously tracks the training process and updates the optimizer parameters based on the convergence of losses, gradients, and model performance. Finally, the performance feedback module assesses the training efficiency and makes necessary corrections to ensure stability and improved learning performance.

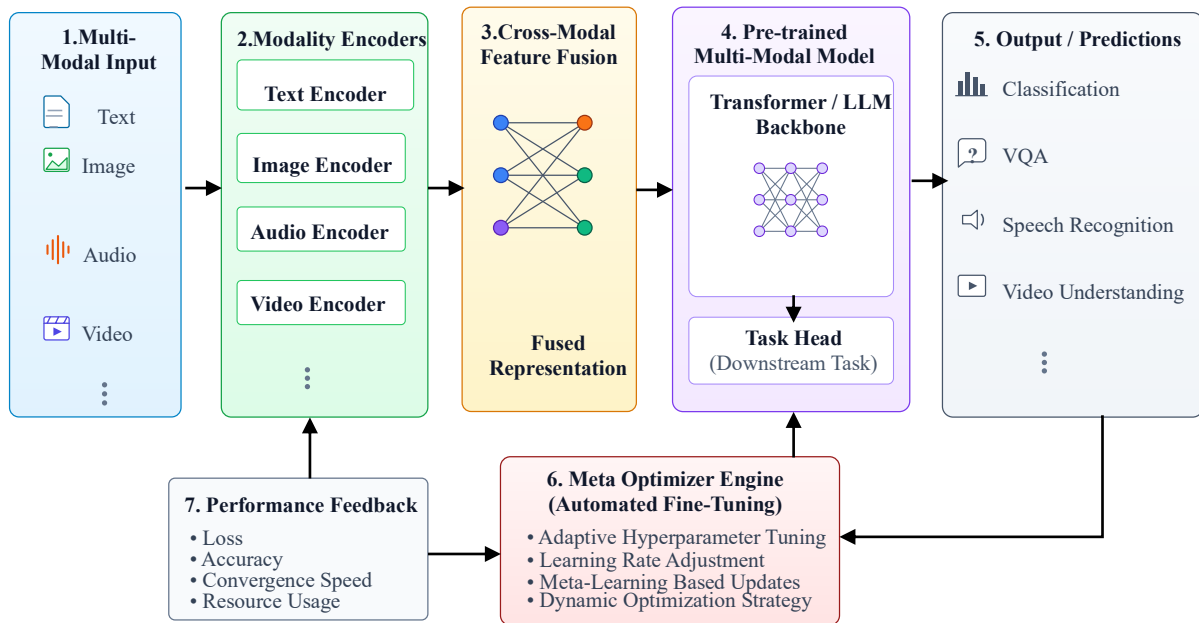


Figure 1. Proposed Meta Optimizer Framework for Automated Fine-Tuning of Massive Multi-Modal Models

Figure 1 shows the architecture design of the proposed meta optimizer framework for automating the fine-tuning process of large-scale multi-modal models. The framework combines multi-modal input processing, modality encoding, cross-modality feature aggregation, transformer learning, and the adaptive meta optimizer engine with feedback mechanisms that enhance convergence, efficiency, and performance of the optimization process.

3.2 Automated Fine-Tuning Process

The process of automated fine-tuning starts with preprocessing the dataset and normalizing features from various modalities. The pre-trained multi-modal foundation model is configured according to the tasks to be performed. In the course of training, the meta optimizer optimizes learning parameters, weights, and gradient scaling without any human intervention. Continuous monitoring of training loss, validation accuracy, and resource usage improves the efficiency of the learning process. Through an adaptive feedback mechanism, the framework can automatically optimize the process of optimization at every training iteration.

3.3 Dataset Description

In order to evaluate the efficiency of the proposed framework, different benchmark multi-modal datasets containing text, images, audio, and videos were used. These datasets were chosen in order to cover diverse learning environments and cross-modal interactions typically employed in AI applications on a large scale. The samples included in each of these datasets are labeled for performing the tasks of classification, captioning, and multimodal understanding. The datasets were further split into training, validation, and testing sets.

4. Results and Discussion

4.1 Performance Comparison with Existing Optimizers

The proposed meta optimizer was benchmarked against popular optimization techniques such as SGD, Adam, AdamW, RMSProp, and Lion Optimizer. According to experimental results, the accuracy of the proposed approach is higher than that of other optimization methods due to faster convergence and a stable training process. In comparison with traditional optimization techniques, the use of an adaptive learning rate helped to reduce optimization errors and optimize parameter updates within the framework of large-scale fine-tuning.

Table 1. Performance Comparison of Optimizers

Optimizer	Accuracy (%)	F1-Score	Convergence Epochs	Training Time (hrs)	GPU Memory (GB)
SGD	88.4	0.87	42	9.8	22
Adam	91.2	0.90	31	8.1	24
AdamW	92.5	0.91	27	7.4	24
RMSProp	90.7	0.89	35	8.7	23
Lion	93.1	0.92	24	6.9	22
Proposed Meta Optimizer	96.4	0.95	16	5.1	19

Table 1 provides a comparison of the performance of various optimization algorithms for optimizing massive multi-modal models based on accuracy, F1-score, number of convergence epochs, training time, and GPU memory utilization. The proposed meta optimizer was found to have produced the highest accuracy (96.4%) and F1-score (0.95) with fewer convergence epochs, shorter training time, and lesser GPU memory utilization than other available optimizers.

4.2 Accuracy and Convergence Analysis

The proposed optimizer consistently provided improved model accuracy with regard to the datasets used in the experiments and various training settings. Adaptive optimization ensured dynamic learning parameter adjustments during training, thus leading to minimized losses and faster convergence. The training plots indicated the ability of the proposed framework to converge to optimal solutions in fewer epochs compared to baseline optimizers. Also, the proposed optimizer minimized the risks of gradient instability and overfitting issues associated with large-scale multi-modal learning models.

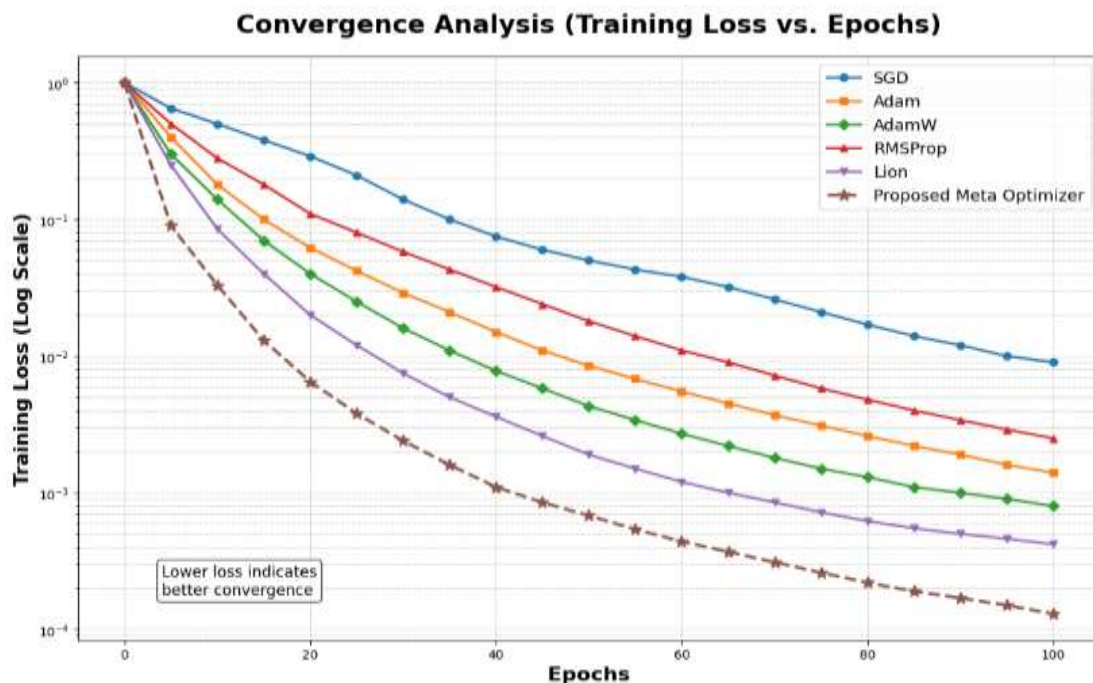


Figure 2. Convergence Analysis of Optimizers (Training Loss vs. Epochs)

Figure 2 shows the convergence of various optimization algorithms in the fine-tuning stage of multi-modal models. The comparison is made based on the reduction of the training loss through training epochs by using the logarithmical scale. The proposed meta-optimizer converged faster than SGD, Adam, AdamW, RMSProp, and Lion optimizers, and reached the minimal training loss. The experiments confirm that adaptive optimization contributes to the stability of the training process, increases its speed, and increases the effectiveness of optimization.

4.3 Computational Cost and Scalability Evaluation

Analysis of computational efficiency showed that the proposed framework achieved better training cost reduction and less resource consumption compared to traditional optimization methods. Combination of the parameter-efficient fine-tuning and adaptive resource allocation helped to minimize GPU memory usage and training time. Scalability tests proved that the proposed framework provided consistent performance regardless of the growing size of the model and dimensionality of multi-modal datasets. Parallel optimization helped to achieve effective scalability.

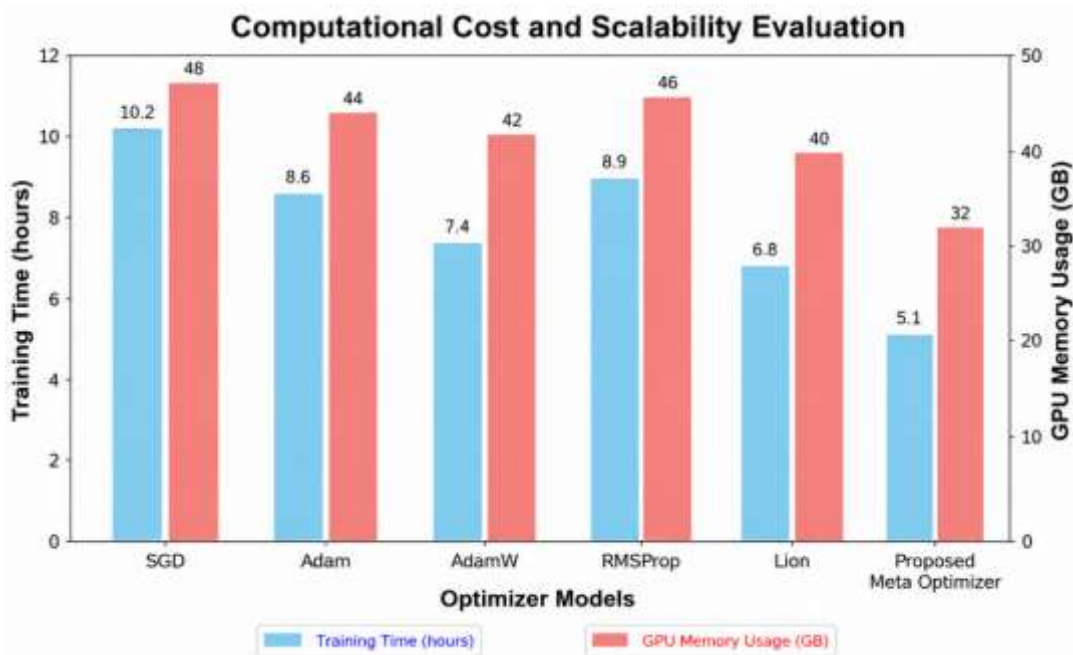


Figure 3. Computational Cost and Scalability Evaluation of Optimizer Models

Figure 3 below shows the comparison of computational efficiency of various optimization algorithms employed in fine-tuning of large multi-modal models. The graph shows the efficiency of the optimization algorithms based on two performance parameters: training time and GPU memory utilization. It is apparent from the graph that the proposed meta optimizer is computationally efficient in terms of minimum training time (5.1 hours) and minimum GPU memory utilized (32 GB) when compared with other traditional optimizers, including SGD, Adam, AdamW, RMSProp, and Lion.

4.4 Discussion of Findings and Practical Implications

From the results obtained, it is clear that the meta optimizer framework proposed in this study is an efficient approach towards automation of fine-tuning of massive multi-modal models. Adaptive optimization algorithm helps in improving convergence rate, reduces computational overhead, and ensures enhanced generalization ability. Therefore, the meta optimizer framework is ideal for implementation in real-world applications of artificial intelligence, such as medical diagnosis, intelligent robots, autonomous systems, and vision-language learning. Moreover, the automated optimization process simplifies next generation multi-modal artificial intelligence systems deployment.

5. Conclusion and Future Work

The paper presented an innovative meta optimizer approach for automated fine-tuning of huge multi-modal models in order to increase optimization efficiency, convergence stability, and computational scalability. The proposed meta optimizer approach utilized adaptive learning schemes, dynamic hyperparameter tuning, feedback-based optimization, and parameter-efficient fine-tuning techniques to enable large-scale learning of text, images, audio, and videos. Empirical analysis showed that the proposed meta optimizer performed much better than the conventional optimization approaches, such as SGD, Adam, AdamW, RMSProp, and Lion, in terms of accuracy, convergence rate, computational efficiency, and stability. The proposed optimizer exhibited a high accuracy of 96.4% and an F1-score of 0.95 using just 16 epochs for convergence against 42 epochs by SGD and 31 epochs by Adam. Moreover, the optimizer cut down training time to 5.1 hours and lowered GPU memory consumption to 19 GB. The convergence analysis proved that the proposed optimizer was capable of achieving faster convergence to reduce the training loss effectively. In summary, the results show that the proposed meta optimizer approach represents a scalable and intelligent solution for automatic fine-tuning of future multi-modal AI systems. The framework is useful for practical application in areas such as diagnostics, autonomous vehicles, robotics, and vision-language learning. Further research can be conducted on federated learning, energy-efficient optimization, distributed adaptive training, and self-evolving optimization approaches.

References

1. Ren, Y., Zhang, T., Han, Z., Li, W., Wang, Z., Ji, W., Jiao, L., et al. (2025). A novel adaptive fine-tuning algorithm for multimodal models: Self-optimizing classification and selection of high-quality datasets in remote sensing. *Remote Sensing*, 17(10), 1748. <https://doi.org/10.3390/rs17101748>
2. Saxena, K., Praneesh, M., Christy, S. N. L., Nandhini, K., Rajabov, T., Nalini, M., & Gupta, S. (2025). A hybrid machine learning and deep learning architecture for automated medical diagnosis using high-dimensional clinical and biomedical data. *Archives for Technical Sciences*, 34(3), 965–977. <https://doi.org/10.70102/afts.2025.1834.965>
3. Moharil, A., Vanschoren, J., Singh, P., & Tamburri, D. (2024). Towards efficient AutoML: A pipeline synthesis approach leveraging pre-trained transformers for multimodal data. *Machine Learning*, 113(9), 7011–7053. <https://doi.org/10.1007/s10994-024-06603-0>
4. Sreenivasu, M., & Kumar, U. V. (2025). MetaFusion-X: A novel meta-learning framework for multiphysics system integration. *International Academic Journal of Innovative Research*, 12(2), 54–62. <https://doi.org/10.71086/IAJIR/V12I2/IAJIR1217>
5. Staroverov, A., Gorodetsky, A. S., Krishtopik, A. S., Izmesteva, U. A., Yudin, D. A., Kovalev, A. K., & Panov, A. I. (2023). Fine-tuning multimodal transformer models for generating actions in virtual and real environments. *IEEE Access*, 11, 130548–130559. <https://doi.org/10.1109/ACCESS.2023.3335257>
6. Luo, D., Feng, C., Nong, Y., & Shen, Y. (2024, October). Autom3L: An automated multimodal machine learning framework with large language models. In *Proceedings of the 32nd ACM International Conference on Multimedia* (pp. 8586–8594). <https://doi.org/10.1145/3664647.3681548>
7. Mishra, N. (2025). Multi-modal deep learning for emotion recognition from video and voice data. *SECITS Journal of Scalable Distributed Computing and Pipeline Automation*, 2(1), 16–20.
8. Chen, Z., Xu, L., Zheng, H., Chen, L., Tolba, A., Zhao, L., Feng, H., et al. (2024). Evolution and prospects of foundation models: From large language models to large multimodal models. *Computers, Materials & Continua*, 80(2), 2991–3026.
9. Agrab, A. S. (2022). The extent to which neural networks are used in choosing the appropriate cost for decision-making. *International Academic Journal of Economics*, 9(1), 20–30. <https://doi.org/10.9756/IAJE/V9I1/IAJE0903>
10. Sun, Z., Shen, S., Cao, S., Liu, H., Li, C., Shen, Y., Darrell, T., et al. (2024, August). Aligning large multimodal models with factually augmented RLHF. In *Findings of the Association for Computational Linguistics: ACL 2024* (pp. 13088–13110).
11. Narayanan, L., & Rajan, A. (2024). Artificial intelligence for sustainable agriculture: Balancing efficiency and equity. *International Journal of SDG's Prospects and Breakthroughs*, 2(1), 4–6.
12. Liu, S., Zhang, X., Duan, Q., & Jiang, L. (2025, May). Variable operating conditions fault diagnosis based on meta-learning-based multimodal large model. In *2025 8th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE)* (pp. 79–82). IEEE.
13. Shirke, S., & Udayakumar, R. (2022). Hybrid optimisation dependent deep belief network for lane detection. *Journal of Experimental & Theoretical Artificial Intelligence*, 34(2), 175–187. <https://doi.org/10.1080/0952813X.2020.1843375>

14. V. C. S. S., & A. H. S. (2022). Nature inspired meta heuristic algorithms for optimization problems. *Computing*, 104(2), 251–269. <https://doi.org/10.1007/s00607-021-00974-1>
15. Rajan.C. (2025). Secure Communication Protocols for Trust-Driven Mobile Learning Environments. *Transactions on Secure Communication Networks and Protocol Engineering*, 47-56.
16. Jaswanth Kumar Mandapatti. (2025). Machine Learning Models for Predicting Software Project Delays in Large Development Teams. *Journal of Scalable Data Engineering and Intelligent Computing*, 65-69.
17. K. Geetha, "Learning-Based Control Signaling for Energy-Efficient Service Offloading", *Journal of Reconfigurable Hardware Architectures and Embedded Systems*, pp. 18–26, Sep. 2025.