



Challenges of Vision-based Control in Robotic Application Utilising Computer

¹Nurudeen Danladi Garba*, ²Muhammad Muneer Haruna, ³Victor Julius, ⁴Samson Ibukun Ogundeji, ⁵Ikenna Ukpai Iro, Aliyu Surajo⁶

^{1,2,3,4,5,6} Department of Mechatronics Engineering, Nigerian Defence Academy (NDA), Kaduna, Nigeria.

DOI: 10.5281/zenodo.15625534

Submission Date: 30 April 2025 | Published Date: 09 June 2025

*Corresponding author: **Nurudeen Danladi Garba**

Department of Mechatronics Engineering, Nigerian Defence Academy (NDA), Kaduna, Nigeria.

Abstract

Artificial intelligence (AI) techniques including deep learning, transformers, deep reinforcement learning (DRL), and large language models (LLMs) have greatly improved robotic capabilities in recent years. Convolutional Neural Networks (CNNs), Vision Transformers (ViTs), DETection Transformers (DETR), the YOLO family of algorithms, segmentation methods, and 3D vision technologies are some of the major AI models propelling advances in robotic vision. Robots can now comprehend and produce language similar to that of humans thanks to LLMs, which use enormous volumes of text data to improve robot-human interaction. This improves the robots' capacity for cooperation and communication in a variety of applications. Robots can now understand and produce human-like language because to LLMs, which use vast volumes of text data to improve human-robot interactions. This improves the robots' ability to communicate and work together in a variety of applications. By improving robotic systems' overall performance, efficiency, and flexibility, the integration of these AI techniques opens the door for increasingly complex and intelligent autonomous agents. Robots and other mechatronic systems can interact with their surroundings and carry out tasks more precisely and flexibly thanks to vision-based control, which employs computer vision to direct their movements. Continuous measurement and feedback are made possible by this method, which increases the system's resilience to mistakes and its ability to adjust to environmental changes.

Keywords: Artificial intelligence (AI), Computer Vision, Mechatronics Systems, Machine Learning (ML), Robotics, Deep Learning (DL).

1. INTRODUCTION

Over the past five years, deep learning (DL) models—especially those in computer vision (CV)—have improved quickly. With the use of deep learning (DL) techniques, CV, a type of artificial intelligence (AI), allows machines to identify and comprehend images. With the aid of data, computer processing power, and machine learning advancements, artificial intelligence (AI) has been steadily advancing and becoming more effective globally, particularly during the past 20 years. It is therefore not surprising that artificial intelligence (AI) has many applications in the military sector as well, in a wide range of fields. AI is being used more and more in the daily lives of many different sectors, including speech recognition, biometric authentication, mobile mapping, navigational systems, transportation and traffic control, management, manufacturing, supply chain management, data collection, and targeted online marketing [1].

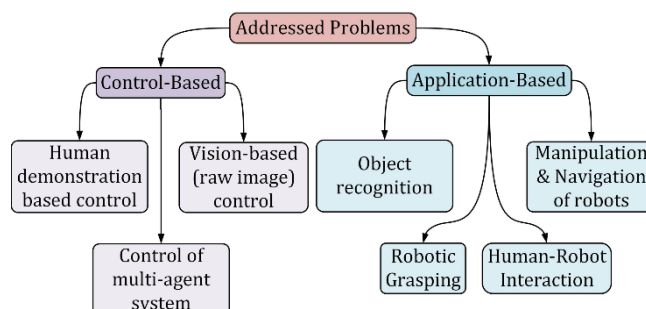


Fig. 1: A Model of Vision-Based Robotic Applications System.

II. LITERATURE REVIEW

Associated work Theoretical underpinnings and creative applications of vision-based robot control have grown significantly in recent years due to developments in deep learning (DL) and camera technologies. A few of the most recent advancements and their uses are highlighted in this succinct assessment, especially as they relate to visual information-based robot control. Picture pre-processing, feature extraction, classification, and outcomes analysis are the four main steps in the image classification process. Raw images should be enhanced during the pre-processing phase of image processing in order to increase feature extraction accuracy and image classification outcomes. Gathering pertinent information from an image in order to construct distinct classes is the main goal of feature extraction [9]. As a result, deep learning-based approaches are both a good substitute and the current standard [10]. Nonetheless, recent studies on convolutional neural networks (CNN) have demonstrated superior generic visual descriptions that are applicable to a variety of picture identification uses, such as the classification of fine-grain images [11]. The methodology of this article is first focused on identifying a good computer vision-based technique for the classification of sensitive documents from non-sensitive documents. The authors first identified the standard parameters in terms of reliability, loss, precision, and recall by using deep learning techniques, such as neural networks with convolutions and transfer learning (TL) algorithms. Many security concerns require robotic vision automation for the segregation of sensitive and non-sensitive documents. Numerous publications have cited the extraction of features based on pre-trained deep learning models. Similarly, we used the majority of feature extraction techniques to identify features from the images, which were subsequently classified by machine and ensemble learning models. Nevertheless, the feature extraction and machine learning classification based on pre-trained models produced better results than the deep learning and TL procedures. Additionally, the more well-known methods were used as the brains behind a robotic structure's vision to automatically separate sensitive papers from non-sensitive documents. When we need to locate a specific, classified document in the haystack, we could use the suggested robotic framework [12]. A 3D stereo vision system for real-time depth perception, a deep reinforcement learning (RL) model tailored for suturing tasks, and a robotic arm with a needle driver make up the suggested framework. In our method, a reinforcement learning agent is trained in a simulated environment to learn the best suturing techniques through trial-and-error interactions. In order to maximize precision and minimize tissue injury, the RL model takes into account tissue deformation, suture tension, and needle trajectory. Accurate needle insertion and suture placement are made possible by the 3D vision module's high-resolution depth maps, which are used to guide the robot in real time. When tested on artificial tissue models, the technology performs better in terms of accuracy, uniformity of sutures, and flexibility in response to tissue changes. According to experimental results, our RL-based method works better than conventional teleoperated suturing by lowering variability and obtaining higher accuracy. This study demonstrates the viability of autonomous robotic suturing in spite of obstacles including dynamic tissue behavior and real-time processing limitations. Future developments will concentrate on increasing the system's adaptability in real time, maximizing computational efficiency, and broadening its scope of use to include more surgical procedures. An important step toward completely autonomous robotic surgery is represented by this study [13]. Significant progress has been made in automated surgical suturing, which combines robots, artificial intelligence, and optimization methods to improve accuracy, effectiveness, and safety. Early studies concentrated on creating mechanical suturing devices that outperformed manual suturing techniques in terms of precision and consistency [14]. Additionally, autonomous flexible endoscopes were developed to enable safer and more flexible minimally invasive operations. Robots were trained to conduct suturing using deep reinforcement learning and pre-learned motion trajectories in order to mimic human suturing abilities using trajectory learning and non-rigid registration techniques [15]. The precision and effectiveness of multi-throw suturing were further enhanced by optimization methods such sequential convex optimization and mechanical needle guides [16]. In order to guarantee safety, accountability, and compliance, ethical, legal, and regulatory considerations become increasingly important as surgical automation advances. Surgical robots has also included machine learning techniques, which allow for better dexterity in autonomous suturing, improved visual perception, and real-time decision-making [17]. More accurate and controlled robotic treatments are now possible thanks to advancements in endoscopic tissue manipulation and needle tracking made possible by vision-based guidance. Even though these developments have a lot of potential, there are still issues with guaranteeing flexibility, legal compliance, and prompt response in intricate surgical settings [18]. To prepare the way for completely autonomous robotic surgery, future research should concentrate on enhancing real-time perception, perfecting AI-driven techniques, and guaranteeing ethical surgical automation. In this study [19], we first describe the continuous creation of a flexible robotic cell that is made possible by deep-learning vision. This cell will be used to automate three primary functions in circular healthcare: the mapping and quantification of resources, the disassembly of minor medical devices, and the sorting of waste. Second, we incorporate robotics into a system-level perspective by fusing the mechanics of robots with compartmental dynamical thermodynamics. By utilizing dynamical systems theory and adding dynamical energy balances to the standard mass balances, our thermodynamic framework improves material flow analysis (MFA) and advances the theoretical underpinnings of circular material flow designs. Third, using graph theory and our thermodynamic framework, we provide two circularity indicators. The robotic cell can be used for reprocessing glucose meters and inhalers, but with the right modifications, it may also be used for other medical devices. It can also be used for sorting, disassembly, resource mapping, and quantification, or it can be used in simultaneously. At the expense of additional complexity, our thermodynamic systemic modeling framework can produce the necessary increases in model accuracy

and reproducibility since it incorporates more physics and system dynamics than MFA. Lastly, healthcare chain managers can use the suggested circularity indicators to determine whether the robotic cell can process the input stream of materials in the allotted time and with the appropriate degree of separation at the output material flow. A demo video and software are made available to the general audience. In addition to discussing the associated difficulties and potential future directions for the field of computer vision in enhanced robotic perception, this review paper provides an overview of the basic ideas of CNNs and their applications in various computer vision tasks for robotic perception. The history, fundamental ideas, operation, uses, and key elements of CNNs are all covered in this essay. To maximize CNNs' potential in robotic perception and enhance robotic performance, it is crucial to comprehend their principles, benefits, and drawbacks [20]. Early autonomous robot movement prediction can help with command issuance to monitor and control robots' future actions before they happen, according to this study [2]. In order to do this, we suggest a deep learning classifier and effective computer vision approaches. The computer vision techniques are intended for feature estimation, object segmentation, and frame quality enhancement. Robot motions are automatically detected by the Long-Term Short Memory (LSTM) classifier using initial sequential features. In order to overcome the issues of vanishing and exploding gradients and to execute the earlier prediction from the initial sequential features of partial video frames, we primarily used an LSTM classifier in the design of the suggested model. Higher accuracy prediction time reduction is facilitated by LSTM. Additionally, it allows a particular robotic application's central system to avoid collisions brought on by obstacles in an indoor or outdoor environment. We sought to develop deep learning (DL) computer vision artificial intelligence (AI) models in [3] that could automatically evaluate trainee performance and determine robotic suturing task competency. The basic ideas of convolutional neural networks (CNNs) and their uses in various computer vision tasks for robotic perception are summarized in [4], along with the associated difficulties and potential directions for the field of computer vision in enhanced robotic perception. The history, fundamental ideas, operation, applications, and key elements of CNNs are also covered. To fully utilize CNNs' potential in robotic perception and improve robotic performance and intelligence, it is essential to comprehend their principles, advantages, and limitations. By removing the complications associated with previous systems, such as laborious mechanics, rigidity, the need for several sensors, etc., this study [5] aims to build an object-tracking system that is both precise and responsive. With vision-based control, the robotic arms can be efficiently built to track moving objects on their own. Image-based visual servoing (IBVS) is more beneficial in vision-based control when compared to other classical and traditional servoing techniques. This article outlines a novel method for IBVS-based tracking control of a robotic arm with two degrees of freedom (DOF) that incorporates trajectory tracking and object identification as essential elements. An accurate deep learning-based object detection framework is used to address the problems related to IBVS. The objects are detected and located in real-time using the framework that has been described. Furthermore, using the object detection system's real-time response, an efficient vision-based control method is created to operate the 2-DOF robotic arm. Using the CoppeliaSim robot simulator and a 2-DOF robotic arm, respectively, modeling and experimental studies are conducted to validate the suggested control technique. The results show that when executing visual servoing tasks, the suggested deep learning controller for the vision-based 2-DOF robotic arm achieves good levels of accuracy and response time [5]. We built an autonomous driving system that can avoid obstacles and drive without the use of a global navigation satellite system (GNSS) signal, which is presented in [6]. Because deep learning techniques, such as convolutional neural networks, can be used to understand an environment and avoid obstacles, the system uses an object identification system based on a stereo camera. Through studies carried out in the University of Tokyo's Tanashi Forest, the vehicle's autonomous driving capabilities were assessed utilizing real-time kinematic-GNSS to measure the genuine values. In order to improve robotic manipulation, this study introduces an intelligent robotic object grasping system that makes use of deep reinforcement learning and computer vision techniques. You Only Look Once (YOLOv3) is a real-time object detection and localization technology presented by the authors in [7]. The Soft Actor-Critic (SAC) system leverages depth image information to identify the best gripping areas. The robotic manipulator may then effectively choose and arrange things by converting the gripping point into a three-dimensional grasping position. YOLO's detection accuracy was improved by using the COCO dataset, and the training process was accelerated by transfer learning. The suggested system's performance evaluation showed an 87.3% grasping success rate and a mean Average Precision (mAP) of 91.2% for item detection. The model's robustness and generalizability were confirmed by 10-fold cross-validation, which showed little difference in performance across test circumstances. The suggested method increased execution efficiency by 35% and accuracy by 27% when compared to conventional gripping techniques. These results show the potential of the YOLO-SAC framework for real-world robotic applications by offering a scalable and adaptable method for automated object handling in many contexts. In this study [8], we introduce an eye-in-hand camera-based vision-based pick-and-place control system for industrial robots. Robots equipped with cameras significantly increase productivity and performance in the workplace. The employment of robotic arms for pick-and-place tasks in simulated settings has been the subject of earlier research. Aligning the coordinate systems between the robot and the camera and maintaining high data accuracy throughout the experiment are the challenges that come with working with real systems. Our research focuses on using deep learning algorithms mounted on the robotic arm's end-effector in conjunction with a low-cost 2D camera to overcome this problem. Both simulation and real-world experiments are used to assess this study. We propose a novel approach that combines the YOLOv7 (You Only Look Once V7) deep learning network with GAN (Generative Adversarial Networks) to achieve fast and accurate object recognition. This system uses deep learning to process camera

data to extract object positions for the robot in real-time. Due to its advantages of fast inference and high accuracy, YOLO is applied as the baseline for research. By training the deep learning model on diverse objects, it effectively recognizes and detects any object in the robot's workspace. Through experimental results, we demonstrate the feasibility and effectiveness of our vision-based pick-and-place system. Our research contributes an important advancement in the field of industrial robots by showcasing the potential of using a 2D camera and an integrated deep learning system for object manipulation.

III. KEY APPLICATIONS OF VISION-BASED CONTROL SYSTEMS

1. Robotics:

Robotic Arm Control: Vision-based control enables robotic arms to perform tasks like assembly, packaging, and inspection with greater accuracy and efficiency.

Visual Servoing: This technique uses feedback from a camera to control the motion of a robot, allowing it to perform fine positioning tasks or track moving objects.

Mobile Robot Control: Vision can be used to guide mobile robots along pre-taught paths or navigate in real-time based on their visual perception.

2. Industrial Automation:

Quality Control: Vision systems can inspect products on a production line to ensure they meet quality standards, detect defects, or identify variations.

Object Recognition and Sorting: Vision can be used to identify and sort objects based on their features, colors, or other visual characteristics.

Barcode Reading: Vision systems can efficiently read barcodes on packaging, labels, or other objects.

3. Other Applications:

Medical Field: Vision-based control can be used in surgery, wound care, and rehabilitation.

Agriculture: Vision can be used for tasks like planting, harvesting, and monitoring crop growth.

Autonomous Vehicles: Vision systems are crucial for autonomous vehicles to perceive their environment, detect objects, and make driving decisions.

Space Applications: Vision-based control can be used for tasks like spacecraft navigation, object tracking, and remote sensing [21].

IV. CHALLENGES OF VISION BASED CONTROL IN ROBOTIC APPLICATION UTILISING COMPUTER

Vision-based control in robotics, utilizing computers for image processing, faces several challenges. These include the inherent variability and complexity of real-world environments, limitations in contextual understanding by computer vision systems, and the need for robust algorithms that can generalize to unseen scenarios. Data quality and annotation, high computational costs, and real-time processing limitations also pose significant hurdles.

Here's a more detailed breakdown of the challenges:

1. Environmental Variability and Complexity:

Dynamic lighting:

Robots need to perform reliably even in fluctuating or poor lighting conditions, which can significantly impact image quality and object recognition.

Occlusions:

Objects being partially blocked by other objects or the environment can hinder accurate perception.

Novel objects and cluttered backgrounds:

Robots need to be able to handle a wide range of objects and environments, including those they haven't seen before, which requires robust algorithms that can generalize.

Real-world variations:

Robots must adapt to a multitude of real-world scenarios, including variations in textures, materials, and the presence of unexpected objects.

2. Limitations in Contextual Understanding:

Lack of semantic understanding:

While computer vision can identify objects, it often struggles to understand the context and relationships between objects in a scene.

Limited predictive reasoning:

Robots need to be able to not only recognize objects but also predict their behavior and anticipate their impact on the environment.

3. Data and Computational Challenges:

Data quality and annotation:

Training computer vision models requires large amounts of high-quality, annotated data, which can be difficult and time-consuming to acquire.

Computational cost:

Real-time processing of complex visual information can be computationally expensive, requiring powerful hardware and efficient algorithms.

Real-time processing limitations:

The need for real-time processing puts pressure on algorithm speed and hardware performance.

Generalization across tasks:

Algorithms trained for one task may not generalize well to other tasks or environments, requiring extensive retraining or adaptation.

4. Other Challenges:

Calibration issues:

Ensuring accurate camera calibration for different robots and environments can be a challenge.

Human-robot interaction:

Robots need to be able to interact safely and effectively with humans, which requires understanding human intentions and adapting to their actions.

Safety and reliability:

Ensuring the safety and reliability of vision-based robotic systems in unpredictable environments is paramount.

Overcoming these challenges requires advancements in computer vision algorithms, hardware, and control systems, as well as ongoing research into contextual understanding and human-robot interaction [22].

V. CONCLUSION

In this study, we have examined numerous articles about vision-based control systems employing a variety of methods, as well as machine learning and deep learning. We have also spoken about how they have advanced their technology through various technological approaches, and we have also talked about their future plans. The difficulties that vision-based control faces in computer-based robotic applications are also discussed.

VI. REFERENCES

1. Das, S., Dey, A., Pal, A., Roy, N., "Applications of artificial intelligence in machine learning: review and prospect," International Journal of Computer Application, vol. 115, no. 9, pp. 31–41, 2015.
2. Mahajan, H. B., Uke, N., Pise, P. et al., "Automatic robot Manoeuvres detection using computer vision and deep learning techniques: a perspective of internet of robotics things (IoRT)", Multimed Tools Appl 82, pp. 23251–23276, 2023, <https://doi.org/10.1007/s11042-022-14253-5>.
3. Choksi, S., Narasimhan, S., Ballo, M., Turkcan, M., Hu, Y., Zang, C., Farrell, A., King, B., Nussbaum, J., Reisner, A., Kostic, Z., Taffurelli, G., Filicori, F., "Automatic assessment of robotic suturing utilizing computer vision in a dry-lab simulation, Art Int Surg, 2025, <https://dx.doi.org/10.20517/ais.2024.84>.
4. Raj, R., Kos, A., "An Extensive Study of Convolutional Neural Networks: Applications in Computer Vision for Improved Robotics Perceptions", Sensors (mdpi), 2025, 25, 1033. <https://doi.org/10.3390/s25041033>.
5. Umesh, K., S., et al., "Autonomous object tracking with vision-based control using a 2DOF robotic arm", 2025, 15:13404, <https://doi.org/10.1038/s41598-025-97930-3>.
6. Kosuke, I., Yutaka, K., Sho, I., Kenji, I., "The Development of Autonomous Navigation and Obstacle Avoidance for a Robotic Mower using Machine Vision Technique", ScienceDirect Science-Direct IFAC, 2019 pp. 173–177,
7. Osita, M., N., Ogochukwu, C., O., Ike, J., M., "Intelligent robotic object grasping system using computer vision and deep reinforcement learning techniques" International Journal of Science and Research Archive, 2025, Vol. 14, No. 03, pp. 511–521, <https://doi.org/10.30574/ijrsra.2025.14.3.0693>.
8. Van-Truong Nguyen, et al., "Vision-Based Pick and Place Control System for Industrial Robots Using an Eye-in-Hand Camera", IEEE ACCESS.
9. Khandan, N., "An intelligent hybrid model for identity document classification," arXiv preprint arXiv:2106.0434, 2021, doi: 10.48550/ arXiv.2106.04345.
10. W. Xiong, X. Jia, D. Yang, M. Ali, L. Li, and S. Wang, "DP-LinkNet: A convolutional network for historical document image binarization," KSII Trans. Internet Inf. Syst. (TIIS), vol. 15, no. 1, pp. 1778–1797, 2021.
11. R. Sicre, A. M. Awal, and T. Furon, "Identity documents classification as an image classification problem," Image Anal. Process. - ICIAP, vol. 2017, pp. 602–613, 2017.
12. Vikas, K., et al., "Deep trained features extraction and dense layer classification of sensitive and normal documents for robotic vision-based segregation", DE GRUYTER, 2024, <https://doi.org/10.1515/pjbr-2022-0125>.
13. Chetan, G., et al., "Autonomous Suturing in Robotic Surgery Using Reinforcement Learning and 3D Visual Feedback", Journal of Neonatal Surgery, Vol. 14, No. 10, pp. 24–35, 2025.
14. A. Shademan, R. S. Decker, J. D. Opfermann, S. Leonard, A. Krieger, and P. C. W. Kim, "Supervised autonomous robotic soft tissue surgery," Sci. Transl. Med., vol. 8, no. 337, pp. 337–364, May 2016.
15. K. Watanabe, T. Kanno, K. Ito, and K. Kawashima, "Single-master dual-slave surgical robot with automated relay of suture needle," IEEE Trans. Ind. Electron., vol. 65, no. 8, pp. 6335–6343, Apr. 2018.
16. Li Y, Richter F, Lu J, et al. "SuPer: A surgical perception framework for endoscopic tissue manipulation with surgical robotics", IEEE Robot Autom Lett. 2020, Vol. 5, No. 2, pp. 2294–2301.

17. S. Petscharnig and K. Schöffmann, “Learning laparoscopic video shot classification for gynecological surgery,” *Multimedia Tools Appl.*, vol. 77, no. 7, pp. 8061–8079, Apr. 2018.
18. S. Speidel, A. Kroehnert, S. Bodenstedt, H. Kenngott, B. Muller-Stich and R. Dillmann, “Image-based tracking of the suturing needle during Laparoscopic interventions,” *Proc. SPIE Med. Imag.*, vol. 9415, Apr. 2015, Art. no. 94150B.
19. Federico, Z., Denis, S., Shahin, R., “Towards a Thermodynamical Deep-Learning-Vision-Based Flexible Robotic Cell for Circular Healthcare”, Springer, Circular Economy and Sustainability, <https://doi.org/10.1007/s43615-025-00532-4>.
20. Raj, R.; Kos, A. An Extensive Study of Convolutional Neural Networks: Applications in Computer Vision for Improved Robotics Perceptions. *Sensors* 2025, 25, 1033. <https://doi.org/10.3390/s25041033>.
21. https://www.google.com/search?q=what+is+vision+based+control+used+for&oq=&gs_lcrp=EgZjaHJvbWUqCQgCECMYJxjqAjIJCaaQIxnGOoCMgkIARajGCcY6gIyCQgCECMYJxjqAjIJCAMQIxnGOoCMgkIBBAjGCcY6gIyCQgFECMYJxjqAjIJCAYQIxnGOoCMgkIBxajGCcY6gLSAQk1MzAwajBqMTWoAgiwAgHxBb1djv6ADJdO&sourceid=chrome&ie=UTF-8.
22. https://www.google.com/search?q=challenges+of+vision+based+control+in+robotic+application+utilising+computer&sca_esv=a8f2e26b4241c72c&sxsrf=AHTn8zoJSmCH0vJjCpMdl44jv8Rh6as41g%3A1745408015551&ei=D9AlaMG2IYC0hbIPsobU8Qw&oq=challenges+of+vision+based+control+in+robo&gs_lp=Egxnd3Mtd2l6LXNlcniAiKmNoYWxsZW5nZXMGb2YgdmlzaW9uIGJhc2VkIGNvbnRyb2wgaW4gcm9ibyocCCAMyBxAhGKABGAoyBxAhGKABGAoyBRAhGJ8FMgUQIRifBTIFECEYnwUyBRAhGJ8FSNHpA1CTowFY5bsDcAJ4AZABAjgB0wOgAcNRqgEJMi0xOS4xNC4yuAEByAEA-AEBmAIYoALwNsICChAAGLADGNYEGEfCAgcQIxiwAhgnwgIFEAAAY7wXCAGgQABiiBBiJBcICCBAAGIAEGKIEwgIKECEYoAEYwwQYCsICCBahGKABGMMewgIFECEYoAGYAwCIBgGQBgiSBwsyLjAuMTAuMTAuMqAHsYYBsgcJMi0xMC4xMC4yuAfiNg&sclient=gws-wiz-serp.

CITATION

Nurudeen D. G., Muneer, H. M., Victor J., Samson I. O., Ikenna U.I., & Aliyu S. (2025). Challenges of Vision-based Control in Robotic Application Utilising Computer. In *Global Journal of Research in Engineering & Computer Sciences* (Vol. 5, Number 3, pp. 96–101).
<https://doi.org/10.5281/zenodo.15625534>