



JOURNAL OF INFORMATION  
SYSTEMS AND TECHNOLOGY

# Journal of Information Systems and Technology

Vol., No. (2025), 1-8

p-ISSN: XXXX e-ISSN: XXXX

Journal homepage: <https://athallahpublishing.com/index.php/jistech>

Research Paper

## Bridging Vision and Understanding: The Central Role of Computer Vision in AI

Nurul Hidayati

Universitas Amikom, Yogyakarta, Indonesia

\*Corresponding author: [nrlhdyati@gmail.com](mailto:nrlhdyati@gmail.com)

---

### ARTICLE INFO

#### Keywords

Artificial Intelligence

Computer Vision

Machine Understanding

### ABSTRACT

Computer vision has become one of the most critical components in the advancement of artificial intelligence, enabling machines not only to perceive but also to interpret the world around them. This paper explores the central role of computer vision in bridging the gap between visual perception and higher-level machine understanding. By integrating deep learning, pattern recognition, and semantic interpretation, computer vision transforms raw visual data into structured knowledge that supports decision-making, reasoning, and autonomous behavior. The discussion highlights recent progress in image recognition, object detection, scene understanding, and multimodal learning, emphasizing how these innovations drive AI toward more human-like cognition. Furthermore, the paper addresses the challenges of scalability, generalization, and ethical implications, offering insights into future directions for research and applications.

Copyright © 2025 Authors

This is an open access article under [CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/) license



## Introduction

The rapid evolution of artificial intelligence (AI) has significantly reshaped how machines interact with the world, with computer vision emerging as one of its most transformative fields. Unlike traditional computational approaches that rely heavily on structured and symbolic input, computer vision allows machines to process and interpret unstructured visual data—images and videos—that dominate human perception. This capability brings AI closer to human-like cognition, enabling systems not only to “see” but also to “understand.” At its core, computer vision serves as the bridge between perception and understanding. Through advances in deep learning, convolutional neural networks, and multimodal integration, machines can now perform tasks such as object detection, facial recognition, and scene interpretation with unprecedented accuracy. These capabilities extend far beyond recognition, offering insights into context, semantics, and intent—key elements for higher-level reasoning and autonomous decision-making.

The importance of computer vision is reflected in its wide range of applications, from autonomous vehicles and medical imaging to industrial automation and smart surveillance. Each of these domains illustrates the potential of visual intelligence to not only enhance efficiency but also to transform how humans and machines collaborate. However, the journey toward true understanding remains challenging. Issues such as data bias, generalization across domains, interpretability of models, and ethical considerations highlight the complexity of bridging low-level perception with high-level reasoning. This paper aims to explore the central role of computer vision in advancing AI understanding. By examining recent breakthroughs, current challenges, and emerging research directions, we underscore how computer vision continues to shape the trajectory of artificial intelligence toward more adaptive, reliable, and human-centric systems. It is the goal of computer vision, a subfield of AI, to enable computers and other systems to infer relevant information from digital photos, videos, and other visual inputs, and then act upon or recommend that data. Just as AI lets machines reason, computer vision helps computers perceive, observe, and comprehend.

Similar to human vision, computer vision has its limitations, especially when compared to the human starting point. Human vision has the advantage of being trained over many years to distinguish between objects, determine their distance, detect motion, and detect anomalies in an image. Using cameras, data, and algorithms in place of retinas, optic nerves, and a visual cortex, computer vision educates robots to carry out these tasks in a significantly shorter amount of time. A system designed to check items or monitor a manufacturing asset can evaluate thousands of products or processes every minute, allowing it to fast outpace human abilities by picking up on even the smallest faults or difficulties. As a result of its usefulness in fields as diverse as medicine, entertainment, manufacturing, and autonomous vehicles, computer vision has recently become increasingly mainstream. Visual recognition tasks including image ordering, constraining, and identifying are essential to many of these applications. Convolutional neural networks (CNNs) have recently advanced, exhibiting their power with amazing results in best-in-class image recognition assignments and frameworks. The result is

that in computer vision, convolutional neural networks (CNNs) have become the fundamental building blocks of deep learning computations.

Deep Neural Networks (DNN) are a type of neural network that is commonly used in calculations involving computer vision because of their superior picture identification abilities. Visual sign decoding typically makes use of Convolutional Neural Networks (CNN or ConvNet), a subtype of Deep Neural Networks (DNNs). It's also utilized in the fields of computer vision and natural language processing (NLP) to help structure data. A convolutional neural network can be built from a variety of building elements. In this post, we will briefly go over these building pieces, which comprise convolution layers, pooling layers, and fully linked layers. The author then moves on to discuss Deep Learning and other neural network methods. In addition, the book covers Convolutional Neural Networks, their development, and their applications in numerous sectors, including medicine and engineering. Using image identification and object detection technologies, deep learning is a technique for realizing computer vision. Since its inception, computer vision has undergone rapid evolution thanks to the advent of deep learning, greatly improving the precision with which images can be recognized. In addition, an expert system can mimic and recreate the reasoning and decision-making process carried out by human experts. In a way that was impossible with traditional expert systems, machine learning—specifically deep learning—has made it possible to "acquire the tacit knowledge of experts." Using massive data and multiple measurements of a phenomenon, machine learning "systematises tacit knowledge." In this article, we describe some knowledge-based computer vision techniques that incorporate deep learning.

Computer vision is one of the numerous subfields of research that go under the umbrella of image processing, and it has been receiving a lot of attention. It is both an academic field that studies the "realization of vision using computers and an artificial intelligence (AI) field that enables computers and systems to derive meaningful information from digital images, videos, and other visual data and act and make recommendations (allowing humans to do so, as well) based on that information. It is a field that enables computers and systems to derive meaningful information from digital images, videos, and other visual data. If AI makes it possible for computers to think, then computer vision makes it possible for them to see, understand, and observe. This operates in a manner comparable to that of human eyesight and is anticipated to provide humans with an advantage. As a result, the objective of 'computer vision' is to equip computers with capabilities that are analogous to those of the human eye, or to actualize the concept of 'computer vision'. To be more specific, the objective is to develop software for computers that can perform similarly to or even better than human vision by making use of data derived from still images or videos.

Up until this point, the most common form of image processing that computers have been capable of is known as computer graphics, or CG. Computer graphics

(CG) is used for projecting and displaying 3D objects on a 2D display, but computer vision is used to derive 3D information from 2D picture data. This distinction is what distinguishes computer graphics from computer vision. Despite the fact that they represent two distinct technologies, they are complementary to one another and help to further the development of emerging technologies like virtual reality and augmented reality (AR). In augmented reality (AR), three-dimensional computer graphics (3DCG) are superimposed on top of real-world backgrounds, with additional data being supplied by a computer. It is an example of the combination of computer vision, which observes the real environment, with computer graphics, which represent a made-up world.

## **Method**

This study employs a qualitative and exploratory approach to examine the central role of computer vision in advancing artificial intelligence (AI) understanding. The methodology is structured into three main stages: literature review, comparative analysis, and conceptual synthesis. A comprehensive review of academic journals, conference proceedings, and authoritative sources was conducted to identify the state-of-the-art developments in computer vision and AI. The focus was placed on works related to image recognition, object detection, scene understanding, multimodal learning, and semantic interpretation. The time frame of the literature spans the past decade to capture both foundational methods and the most recent advances driven by deep learning and neural architectures.

Key approaches in computer vision—such as convolutional neural networks (CNNs), transformer-based architectures, and hybrid multimodal systems—were compared in terms of their performance, scalability, and ability to support machine understanding. The comparison considered factors including accuracy, generalization capacity, computational efficiency, and applicability to real-world tasks. Case studies from domains such as autonomous vehicles, medical imaging, and surveillance were incorporated to demonstrate practical relevance. Findings from the literature review and comparative analysis were synthesized to develop a conceptual framework that explains how computer vision bridges visual perception and higher-level AI understanding. This framework highlights the transformation of raw visual data into structured knowledge, emphasizing the role of semantic interpretation, contextual reasoning, and integration with other AI modalities (e.g., natural language processing and decision-making systems). The methodology ensures a systematic exploration of the field, balancing theoretical insights with practical applications. By combining evidence from existing studies and synthesizing it into a coherent narrative, the paper aims to provide a comprehensive understanding of the current landscape and future directions of computer vision as a driver of AI understanding.

## Results and Discussion

### Visual information perceived by machine

Computer vision is a multidisciplinary field that enables machines to interpret and understand visual information from the world, much like humans do. It involves the development of algorithms, models, and systems that enable computers to gain a high-level understanding of images or video. Here are some fundamental concepts and components of computer vision. Images are composed of pixels (picture elements), which are the smallest units of an image. The number of pixels in an image determines its resolution, with higher resolution providing more detail. Before analysis, images often undergo preprocessing steps such as resizing, normalization, and noise reduction. Techniques like convolution are used for filtering to enhance or extract specific features from an image. Detection of edges is a common feature extraction technique to identify boundaries in an image. Corner detection helps identify distinctive points in an image. Extracting meaningful information from regions of interest, such as SIFT or SURF descriptors.

**Deep Learning:** Using neural networks, particularly convolutional neural networks (CNNs), for feature learning and representation. *Neural Networks and Deep Learning:* Convolutional Neural Networks (CNNs) are especially popular in computer vision for tasks like image classification, object detection, and segmentation. Understanding how machines perceive visual information involves grasping these fundamental concepts and techniques. The field is continually evolving, with ongoing research and advancements in technology contributing to the development of more robust and versatile computer vision systems. The role of sensors, cameras, and other data sources is crucial in capturing visual data for computer vision applications. These devices serve as the eyes of the computer system, providing the necessary input for analysis and interpretation. Here's a discussion of their roles.

Cameras capture visual information in the form of images. They come in various types, including digital cameras, webcams, and specialized cameras for specific applications. The resolution of a camera determines the level of detail in the captured images. Higher resolution cameras provide more information but may also require more processing power. For video applications, the frame rate of a camera is crucial. Higher frame rates result in smoother videos and are essential for tasks like motion analysis and object tracking. Cameras equipped with depth sensors, such as Time-of-Flight (ToF) or structured light sensors, can capture depth information along with color. This is valuable for tasks like 3D reconstruction and understanding the spatial arrangement of objects. Infrared sensors capture infrared radiation, which is invisible to the human eye. They find applications in night vision, thermal imaging, and other scenarios where detecting heat or radiation is important.

*LiDAR (Light Detection and Ranging):* LiDAR sensors use laser light to measure distances and create detailed, three-dimensional maps of the environment. LiDAR

is commonly used in autonomous vehicles, robotics, and environmental monitoring. *RGB-D Cameras:* These cameras combine traditional RGB (color) information with depth data. They provide a richer representation of the scene and are widely used in applications requiring both color and depth information. Accelerometers and gyroscopes can provide information about the movement and orientation of a device. This data is useful for tasks such as image stabilization, gesture recognition, and tracking dynamic objects. While not visual sensors, microphones capture audio data, which can complement visual information. Audio-visual fusion is used in applications like speech recognition, lip reading, and context-aware computer vision. These sensors capture information across multiple wavelengths, allowing for the analysis of spectral characteristics. They find applications in agriculture, environmental monitoring, and material identification.

Integrating data from multiple sensors can enhance the overall understanding of a scene. Sensor fusion techniques combine information from different sources, improving accuracy and robustness. *Wireless and IoT Devices:* Cameras and sensors integrated into Internet of Things (IoT) devices contribute to the growing field of edge computing, where data is processed locally before being sent to centralized servers. This is particularly useful for real-time applications and reducing latency. In summary, sensors, cameras, and other data sources play a pivotal role in capturing diverse visual data for computer vision applications. The choice of sensors depends on the specific requirements of the task, such as the need for color information, depth perception, or environmental awareness. The continuous advancement of sensor technologies contributes to the ongoing progress in the field of computer vision.

## **Automobile**

Because it is designed to grasp the driving conditions, including spotting barriers, persons on footpaths, and possible accident ways, computer vision is becoming increasingly important to the automobile industry. This is particularly the case given the increased visibility of individuals driving their own vehicles. More and more companies are looking for innovative ways to put more electric vehicles onto the road, which has led to the gradual introduction of self-driving cars onto the market. The advancement of computer vision technology makes it possible for self-driving vehicles to "see" the earth, and artificial intelligence calculations make the "minds" that assist computer vision in translating the objects that are near the vehicle. The newest generation of self-driving cars are outfitted with a plethora of cameras, each of which can provide a full 360-degree view of the natural environment within a range of several meters. For example, the automaker Tesla equips its vehicles with something on the order of 8 all-encompassing cameras to achieve this goal. In addition to the cameras, a front-facing radar that enables the recognition of different vehicles even in the presence of precipitation or mist, as well

as twelve ultrasonic sensors that can distinguish between hard and soft objects that may be found in the environment, have been included.

A standard personal computer won't be adequate to handle the deluge of data that will be pushed into the vehicle because of the amount of information that will be urged into the vehicle. Because of this, every autonomous vehicle has a locally accessible personal computer equipped with computer vision features developed using artificial intelligence. It is the responsibility of the cameras and sensors to identify and collect protests that take place in natural settings, such as groups of people walking. Quick consideration must be given to the dimensions of the items being driven, including their area, thickness, shape, and depth, so that the rest of the driving framework can make the appropriate decisions. All of these computations are only made possible by the combination of AI and deep neural systems, which gives rise to features such as the recognition of a person walking.

Photon estimates, which are mostly made up of photographs of the universe, provide the foundation of our comprehensive knowledge of the cosmos. Because our universe is so massive, and because the one natural rule that governs our world predicts that the data acquired will be just as large, this paves the way for the possibility of employing computer vision in the field of astronomy. It is not possible for the stargazer or anybody else to physically comprehend this information in its entirety in its entirety in its entirety in its entirety. Because of the capabilities of computer vision, we can understand all of the data in a relatively short amount of time. To put it another way, computer vision is currently being used to locate new planets and large bodies. This technique is applied in applications such as imaging of exoplanets, the grouping of stars and cosmic systems, and other activities that are very similar.

## **Conclusion**

Computer vision stands at the core of artificial intelligence's progression from mere perception to genuine understanding. By enabling machines to process, interpret, and contextualize visual data, computer vision bridges the gap between raw sensory input and higher-order cognitive functions. Advances in deep learning, neural architectures, and multimodal integration have not only enhanced the accuracy of image recognition and object detection but have also paved the way for more complex tasks such as scene interpretation, semantic reasoning, and autonomous decision-making. The analysis presented in this paper underscores that computer vision is more than a technical capability—it is a foundational driver of AI's transformation into systems capable of adaptive and intelligent interaction with the world. Its applications in healthcare, transportation, industry, and security illustrate both the immense potential and the profound societal impact of this technology. However, challenges remain. Issues of scalability, domain generalization, interpretability, and ethical responsibility continue to shape the trajectory of computer vision research and practice. Addressing these concerns will

be critical in ensuring that AI systems not only become more powerful but also more trustworthy, fair, and aligned with human values. In conclusion, computer vision is the essential bridge connecting vision to understanding in machines. Its future lies not only in technical refinement but also in fostering responsible integration with other AI domains, ultimately driving artificial intelligence toward more human-centric and cognitively inspired forms of intelligence.

## References

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 25, 1097–1105. <https://doi.org/10.1145/3065386>
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Redmon, J., & Farhadi, A. (2018). YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*. <https://arxiv.org/abs/1804.02767>
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... Fei-Fei, L. (2015). ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3), 211–252. <https://doi.org/10.1007/s11263-015-0816-y>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 30. <https://arxiv.org/abs/1706.03762>
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., ... Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. *International Conference on Machine Learning (ICML)*, 2048–2057. <https://arxiv.org/abs/1502.03044>
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., & Torralba, A. (2016). Learning deep features for discriminative localization. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2921–2929. <https://doi.org/10.1109/CVPR.2016.319>