

# International Conference on Artificial Intelligence and Cloud Computing (ICAICC-2025)

Conference Proceedings

May 10, 2025 (Virtual)

## The Wisdom of Fusion: In-depth Analysis and Future Outlook of Visual Multimodal Technologies

**Dayi Jin**

PhD in Electrical Engineering Specializing in Artificial Intelligence from Stevens Institute of Technology, USA

The application of multimodal computer vision is rapidly evolving, with advancements in deep learning techniques and algorithmic approaches making significant impacts across a variety of industries. This presentation will focus on the cutting-edge algorithms and technologies driving the integration of multimodal data sources to improve visual recognition and image processing. Specifically, we will explore how combining visual, spatial, and temporal data enhances the performance of object detection, recognition, and image segmentation models, enabling more robust systems for real-world applications.

Recent breakthroughs in deep learning, such as transformers and attention mechanisms, have shown promise in overcoming traditional challenges in computer vision, such as handling occlusion, variability in lighting, and dynamic scene changes. These advances are particularly valuable in industries like autonomous driving, where accurate and real-time visual perception is critical for navigation, and healthcare, where image-based diagnostics can be significantly enhanced by incorporating multimodal data from medical imaging devices.

The presentation will delve into key research efforts that integrate multimodal learning for improved performance, focusing on both academic advancements and industry applications. One such example is a framework combining temporal and spatial modalities for gesture and action recognition, which achieved a classification accuracy of 98%. Additionally, we will discuss the application of deep learning models in real-time visual recognition tasks, with examples from healthcare imaging and autonomous vehicle navigation.

By examining the latest research papers, real-world applications, and emerging technologies, this session aims to provide valuable insights into the future direction of multimodal computer vision. It will also highlight interdisciplinary collaborations that are advancing the integration of AI-powered vision solutions into various industries, promoting new capabilities in visual recognition that have the potential to transform how we interact with the world.