

Multimodal learning has emerged as a promising avenue for advancing tasks like scene recognition and few-shot classification by leveraging information from diverse data modalities. Recent works have addressed various challenges in this domain, including effective modality fusion, data sparsity, and cross-modal representation learning.

A notable approach is the two-stage hybrid fusion strategy proposed in **TFAVC**, which combines feature-level and decision-level fusion to integrate audio and visual data robustly. By employing adaptive weighting mechanisms and a weighted fusion embedding, this method captures both the complementarity and relationships between modalities, achieving enhanced performance in scene recognition tasks. The framework's ability to leverage decision-level consensus further ensures resilience to distribution disparities across modalities.

Another critical challenge is modeling the semantic correlations between different modalities. The **SoundingEarth** framework tackles this by encoding images and audio spectrograms into a shared embedding space using modality-specific ResNet architectures. A batch triplet loss optimizes the embeddings, minimizing distances between paired data while maintaining a margin for non-paired examples. This efficient and scalable approach captures the intricate relationships between modalities, enabling effective representation learning. However, its relatively simple architecture limits its ability to handle highly complex cross-modal relationships.

In scenarios where modality data is incomplete or limited, the **HAVE-Net** framework introduces a novel cross-modal hallucination mechanism using conditional multimodal GANs. By generating pseudo-features for missing modalities, this approach enables robust few-shot classification while addressing data sparsity. HAVE-Net further refines class prototypes through metric learning, ensuring consistency even in low-data regimes. While effective, the reliance on adversarial training can introduce instability, especially in large-scale applications.

For aerial scene recognition, **ADVANCE** highlights the importance of preserving modality-specific knowledge while learning cross-modal representations. By employing knowledge distillation to retain audio source knowledge and explicitly modeling scene-sound event relationships, ADVANCE aligns multimodal representations to enhance classification performance. The multi-task learning paradigm further ensures that both modalities contribute meaningfully to the final decision, though the reliance on paired datasets limits its scalability.

These approaches collectively highlight the diversity of strategies for multimodal learning, from robust fusion mechanisms to cross-modal generation and representation learning.