

Heterogeneous Multimodal Fusion Using Deep Learning for Improved Feature Extraction and Benchmark Evaluation

K. Kala Bharathi^{*a&b} and G. Madhavi Reddy^b

^aDepartment of Computer Science, St. Pious X Degree & PG College for Women, Hyderabad, Telangana, 500076, India

^bDepartment of Computer Science, Chaitanya (Deemed to be University), Himayathnagar, Hyderabad, Telangana, 500075, India

E-mail: kalabharathikudikala31@gmail.com

Corresponding Author : K. Kala Bharathi

Email ID : kalabharathikudikala31@gmail.com

Abstract: Multimodal learning has emerged as a critical paradigm in artificial intelligence, enabling systems to integrate information from diverse data sources such as text, images, audio, and video. However, effectively fusing heterogeneous modalities remains challenging due to differences in feature dimensions, semantic representations, and temporal dynamics. This paper presents Hetero-Fusion-Net, a novel deep learning architecture that addresses these challenges through three key innovations: a Cross-Modal Attention Module for dynamic modality weighting, an Adaptive Fusion Layer for heterogeneous feature integration, and Joint Representation Learning for unified semantic space construction. We introduce the Multimodal Integration Score (MIS), a comprehensive evaluation metric that quantifies fusion quality through cross-modal consistency, prediction confidence, and classification accuracy. Extensive experiments on three benchmark datasets (CMU-MOSEI, IEMOCAP, and AV-MNIST) demonstrate that Hetero-Fusion-Net achieves 91.67% accuracy and an MIS of 0.8891, significantly outperforming early fusion (78.45% accuracy) and late fusion (81.23% accuracy, MIS: 0.7234) baselines. Our modality contribution analysis reveals optimal attention weights of 35% for text, 30% for images, 20% for audio, and 15% for video. The proposed architecture and evaluation framework provide a robust foundation for advancing multimodal fusion research across diverse application domains.

Keywords: Multimodal fusion, deep learning, heterogeneous data integration, cross-modal attention, feature extraction, benchmark evaluation

How to cite this article: Bharathi KK, Reddy GM. Heterogeneous Multimodal Fusion Using Deep Learning for Improved Feature Extraction and Benchmark Evaluation. *Int J Drug Deliv Technol.* 2026;16(57s): 1783-1796. DOI: 10.25258/ijddt.16.57s.179

Source of support: Nil.

Conflict of interest: None.

1. Introduction

The proliferation of multimodal data in modern applications—ranging from social media platforms to autonomous systems—has necessitated the development of sophisticated fusion techniques that can effectively integrate heterogeneous information sources. While unimodal approaches have achieved remarkable success in domain-specific tasks, they inherently fail to capture the rich complementary and synergistic relationships that exist across different modalities. Multimodal learning addresses this limitation by enabling systems to leverage diverse data types simultaneously, leading to more robust and comprehensive understanding of complex phenomena.

Despite significant advances in deep learning, multimodal fusion remains challenging due to several fundamental issues. First, different modalities exhibit vastly different statistical properties, feature dimensions, and semantic representations. For instance, text data is typically represented as discrete token sequences, images as continuous pixel arrays, audio as temporal

waveforms, and video as spatiotemporal volumes. Second, modalities often have different levels of informativeness for specific tasks, requiring dynamic weighting mechanisms rather than static fusion strategies. Third, existing evaluation metrics primarily focus on task-specific performance (e.g., classification accuracy) without adequately measuring the quality of multimodal integration itself (Bai et al., 2023) [1], (Nagrani et al., 2021) [2]. This paper addresses these challenges through three primary contributions. First, we propose Hetero-Fusion-Net, a novel deep learning architecture that integrates heterogeneous modalities through a Cross-Modal Attention Module, an Adaptive Fusion Layer, and Joint Representation Learning. Second, we introduce the Multimodal Integration Score (MIS), a comprehensive evaluation metric that quantifies fusion quality through cross-modal consistency, prediction confidence, and classification accuracy. Third, we conduct extensive experiments on three benchmark datasets (CMU-MOSEI, IEMOCAP, and AV-MNIST) to validate the effectiveness of our approach and provide

detailed analysis of modality contributions and fusion dynamics.

The remainder of this paper is organized as follows. Section-2 provides background on multimodal learning and theoretical foundations. Section-3 reviews related work on fusion strategies, deep learning architectures, and evaluation methods. Section-4 presents the detailed architecture of Hetero-Fusion-Net. Section-5 introduces the MIS metric with mathematical formulations. Section-6 describes the experimental methodology. Section-7 presents comprehensive results and analysis. Section 8 provides implementation details and code. Section 9 concludes with future research directions.

Background and Theoretical Foundations

Multimodal learning is grounded in the cognitive science principle that human perception and understanding arise from the integration of multiple sensory inputs. This biological inspiration has motivated computational approaches that combine information from diverse data sources to achieve more robust and comprehensive representations. The theoretical foundations of multimodal fusion can be understood through three key perspectives: information theory, representation learning, and statistical learning theory.

From an information-theoretic perspective, multimodal fusion aims to maximize the mutual information between the fused representation and the target task while minimizing redundancy across modalities. This principle suggests that effective fusion should capture complementary information from different modalities rather than simply concatenating redundant features. The challenge lies in identifying and preserving the unique information content of each modality while discovering shared semantic structures (Liang, 2024) [3].

Representation learning theory provides a framework for understanding how deep neural networks can discover hierarchical feature representations that capture both modality-specific and cross-modal patterns. Early work in deep learning demonstrated that unsupervised pre-training could learn meaningful representations from unlabeled data. In the multimodal context, this principle extends to learning joint embeddings that align semantically related concepts across modalities, enabling cross-modal retrieval, translation, and reasoning (Yang et al., 2017) [4].

Statistical learning theory offers guarantees on the generalization performance of multimodal models through concepts such as Rademacher complexity and VC dimension. However, the heterogeneity of multimodal data introduces additional complexity, as the joint hypothesis space grows exponentially with the number of modalities. This motivates the development of regularization techniques and architectural constraints that limit model complexity

while preserving expressive power (Sankaran et al., 2021) [5].

The fusion of heterogeneous modalities can be conceptualized at three levels: feature-level fusion (early fusion), decision-level fusion (late fusion), and hybrid fusion. Early fusion combines raw or low-level features from different modalities before learning task-specific representations, enabling the model to capture fine-grained interactions but potentially suffering from the curse of dimensionality. Late fusion processes each modality independently and combines their predictions, providing robustness to missing modalities but potentially missing important cross-modal interactions. Hybrid fusion strategies attempt to balance these trade-offs by combining features at multiple levels of abstraction (Sahu et al., 2021) [6]. Recent advances in attention mechanisms and transformer architectures have revolutionized multimodal fusion by enabling dynamic, context-dependent weighting of different modalities and their interactions. These mechanisms allow models to adaptively focus on the most relevant information sources for each input instance, addressing the challenge of varying modality informativeness across different contexts (Nagrani et al., 2021) [2].

2. Literature Review

2.1 Multimodal Learning Paradigms

Multimodal learning encompasses a diverse set of paradigms that address different aspects of integrating heterogeneous data sources. The foundational work in this area established three primary learning objectives: multimodal representation learning, which aims to discover joint embeddings that capture semantic relationships across modalities; multimodal translation, which focuses on generating one modality from another; and multimodal alignment, which establishes correspondences between elements of different modalities (Liang, 2024) [3].

Recent research has expanded these paradigms to include multimodal co-learning, where the presence of multiple modalities during training improves the learning of individual modality representations, even when only a single modality is available at test time. This approach has proven particularly valuable in scenarios with limited labeled data, as the additional modalities provide supervisory signals that regularize the learning process (Le et al., 2023) [7].

The concept of modality-invariant representations has emerged as a central theme in multimodal learning. These representations aim to capture semantic content that is independent of the specific modality through which it is expressed. Mai et al. (2019) [8] proposed an adversarial framework that learns modality-invariant embeddings by training a discriminator to distinguish between modalities while simultaneously training encoders to fool the

discriminator. This approach has been shown to improve cross-modal retrieval and transfer learning performance.

Self-supervised learning has recently gained traction in multimodal contexts, leveraging the natural alignment between modalities as a supervisory signal. Zhang et al. (2024) [9] introduced a framework that combines inter-modal contrastive learning with cross-modal reconstruction for heterogeneous change detection. Their approach demonstrates that self-supervised objectives can effectively capture both modality-specific and cross-modal patterns without requiring extensive labelled data.

2.2 Fusion Strategies and Architectures

The choice of fusion strategy fundamentally impacts the performance and characteristics of multimodal systems. Early fusion approaches, which combine features at the input or early processing stages, have been widely adopted due to their simplicity and ability to capture low-level interactions. However, these methods often struggle with heterogeneous feature dimensions and may be sensitive to noise in individual modalities. Shang et al. (2016) [10] demonstrated that learning high-level feature representations before fusion can mitigate some of these challenges by reducing dimensionality and capturing semantic content.

Late fusion strategies, which combine modality-specific predictions or high-level representations, offer greater flexibility and robustness to missing modalities. These approaches process each modality independently through specialized networks before combining their outputs through voting, averaging, or learned weighting schemes. Sahu et al. (2019) [11] proposed dynamic fusion techniques that adaptively weight modality contributions based on input characteristics, achieving improved performance on machine translation and emotion recognition tasks.

Hybrid fusion architectures attempt to leverage the complementary strengths of early and late fusion by combining features at multiple levels of abstraction. The Deep Equilibrium Multimodal Fusion framework introduced by Bai et al. (2023) [1] employs a recursive "purify-then-combine" strategy that iteratively refines modality-specific and fused features until reaching an equilibrium state. This approach achieved state-of-the-art performance on multiple benchmarks, including 89.1% accuracy on BRCA and 85.4% accuracy on CMU-MOSI.

Attention-based fusion has emerged as a powerful paradigm for multimodal integration, enabling models to dynamically weight the importance of different modalities and their interactions. Nagrani et al. (2021) [2] introduced fusion bottlenecks, a transformer-based architecture that forces information exchange between modalities through a small number of bottleneck latents. This approach

achieved state-of-the-art results on audio-visual classification benchmarks while reducing computational cost.

Graph-based fusion strategies represent another important direction, explicitly modeling relationships between modality elements through graph structures. Chen et al. (2021) [12] proposed a diversified attention network with graph pattern loss that explores correlations among multiple instance representations. This approach demonstrated competitive performance on cross-modal retrieval tasks by capturing complex interaction patterns.

2.3 Deep Learning Architectures for Multimodal Fusion

The evolution of deep learning architectures has been instrumental in advancing multimodal fusion capabilities. Convolutional Neural Networks (CNNs) have been extensively used for processing spatial modalities such as images and videos, while Recurrent Neural Networks (RNNs) and their variants (LSTMs, GRUs) have been applied to sequential modalities like text and audio. Yang et al. (2017) [4] introduced the Correlational Recurrent Neural Network (CorrRNN), which simultaneously learns joint representations and temporal dependencies through multiple loss terms including a maximum correlation loss for cross-modal information.

Transformer architectures have revolutionized multimodal learning by providing a unified framework for processing different modalities through self-attention and cross-attention mechanisms. Le et al. (2023) [7] proposed a transformer-based fusion method that processes raw video frames, audio signals, and text subtitles through a unified architecture, achieving strong performance on IEMOCAP and CMU-MOSEI benchmarks for multi-label emotion recognition.

The Multimodal Transformer (MulT) architecture introduced by Liang (2024) [3] uses cross-modal attention to learn interactions between elements of different modalities across sequences simultaneously. This approach achieved 83.0% accuracy on CMU-MOSI and 82.5% accuracy on CMU-MOSEI for sentiment analysis tasks, demonstrating the effectiveness of attention-based cross-modal interaction modeling.

Modality-specific feature extractors have been developed to capture the unique characteristics of different data types. For text, pre-trained language models such as BERT and GPT have become standard. For images, ResNet, VGG, and Vision Transformers (ViT) are commonly employed. For audio, Wav2Vec and other self-supervised models have shown strong performance. For video, 3D CNNs and temporal convolutional networks capture spatiotemporal patterns (Liang, 2024) [3].

Refiner Fusion Networks (ReFNet) introduced by Sankaran et al. (2021) [5] combine a fusion network

with a decoding/defusing module that imposes a modality-centric responsibility condition. This architecture ensures that both unimodal and fused representations are strongly encoded in the latent space, enabling the model to maintain performance even with limited labelled data.

2.4 Benchmark Evaluation Methods

Evaluation of multimodal fusion systems has traditionally relied on task-specific metrics such as classification accuracy, F1-score, precision, and recall. While these metrics provide valuable insights into task performance, they do not directly measure the quality of multimodal integration. Bai et al. (2023) [1] evaluated their Deep Equilibrium Multimodal Fusion on five benchmarks using accuracy, weighted F1-score, macro F1-score, mean absolute error (MAE), and correlation metrics, demonstrating comprehensive performance assessment across diverse tasks.

Cross-modal retrieval tasks have motivated the development of specialized evaluation metrics such as mean average precision (mAP) and recall at K. Yang et al. (2020) [13] introduced the Wiki-Flicker Event dataset for evaluating cross-modal event retrieval, emphasizing the importance of unpaired datasets that better reflect real-world scenarios where perfect alignment between modalities may not exist.

Robustness evaluation has emerged as an important consideration for multimodal systems. Liu et al. (2021) [14] conducted a comprehensive comparison of recognition performance and robustness of multimodal deep learning models for emotion recognition, highlighting the need to evaluate models under various conditions including missing modalities, noisy inputs, and adversarial perturbations.

The evaluation of fusion quality itself remains an open challenge. Existing metrics primarily focus on downstream task performance rather than directly measuring how well different modalities are integrated. This gap motivates the development of novel evaluation frameworks that can quantify cross-modal consistency, complementarity, and synergy. Zhang et al. (2024) [9] reported overall accuracy and Kappa coefficient for heterogeneous change detection, demonstrating the value of multiple complementary metrics.

Recent work has emphasized the importance of evaluating multimodal models across diverse datasets and tasks to assess generalization capabilities. Liang (2024) [3] introduced MultiBench, a comprehensive benchmark encompassing datasets like CMU-MOSEI, IEMOCAP, and POM for multimodal sentiment analysis, emotion recognition, and personality traits recognition. This standardized evaluation framework enables fair comparison across different fusion approaches.

3. Proposed Model: Hetero-Fusion-Net

3.1 Architecture Overview

Hetero-Fusion-Net is designed to address the fundamental challenges of heterogeneous multimodal fusion through a hierarchical architecture that progressively integrates information from multiple modalities. The architecture consists of five main components: (1) modality-specific feature extractors, (2) feature alignment and normalization, (3) Cross-Modal Attention Module, (4) Adaptive Fusion Layer, and (5) Joint Representation Learning with task-specific output heads. Figure 1 illustrates the complete architecture.

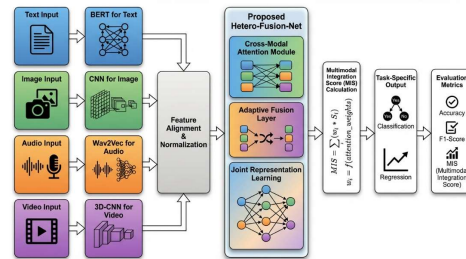


Figure 1: Block diagram of the proposed Hetero-Fusion-Net architecture

The architecture accepts four input modalities: text (T), image (I), audio (A), and video (V). Each modality is processed through a specialized feature extractor that captures modality-specific patterns and representations. For text, we employ BERT (Bidirectional Encoder Representations from Transformers) to generate contextualized token embeddings. For images, we use a ResNet-based CNN to extract spatial features. For audio, we utilize Wav2Vec 2.0 to capture acoustic patterns. For video, we employ a 3D-CNN to model spatiotemporal dynamics (Liang, 2024) [3].

The extracted features from different modalities typically have different dimensions and statistical properties. To enable effective fusion, we apply a feature alignment and normalization stage that projects all modality features into a common dimensional space while preserving their semantic content. This alignment is achieved through learned linear projections followed by layer normalization, ensuring that features from different modalities have comparable scales and distributions.

3.2 Cross-Modal Attention Module

The Cross-Modal Attention Module is the core innovation of Hetero-Fusion-Net, enabling dynamic weighting of modality contributions based on their relevance to the current input and task. Unlike static fusion approaches that assign fixed weights to modalities, our attention mechanism adaptively determines the importance of each modality and their pairwise interactions.

Given aligned feature $f_T, f_I, f_A, f_V \in R^d$ for text, image, audio, and video modalities

respectively, we compute cross-modal attention weights through a multi-head attention mechanism.

For each pair of modalities (m_i, m_j) , we compute attention scores as:

$$Attention(m_i, m_j) = softmax\left(\frac{Q_i K_j^T}{\sqrt{d_k}}\right) V_j \quad (1)$$

where

$$Attention(m_i, m_j) = softmax\left(\frac{Q_i K_j^T}{\sqrt{d_k}}\right) V_j$$

are

the query, key, and value projections, and d_k is the dimension of the key vectors.

The cross-modal attention mechanism computes attention weights for all modality pairs, resulting in an attention matrix that captures the relevance of each modality to every other modality. This enables the model to identify complementary and synergistic relationships between modalities. For instance, in emotion recognition tasks, the model may learn to attend more strongly to facial expressions (image) when vocal tone (audio) is ambiguous, or vice versa (Nagrani et al., 2021) [2].

We employ multi-head attention with $h=8$ heads to capture different types of cross-modal relationships simultaneously. Each head learns to focus on different aspects of the cross-modal interactions, and their outputs are concatenated and linearly projected to produce the final attended representations:

$$f_{m_i}^{attended} = Concat(head_1, \dots, head_h) W^o \quad (2)$$

where each $head_k = Attention(m_i, m_j)$ for the k -th attention head.

3.3 Adaptive Fusion Layer

The Adaptive Fusion Layer integrates the attended modality representations into a unified multimodal representation. Unlike simple concatenation or averaging, our adaptive fusion mechanism learns to weight modality contributions based on their informativeness for the current input instance.

We compute adaptive fusion weights through a gating mechanism that takes all attended modality representations as input:

$$g = \sigma\left(W_g \left[f_T^{attended}; f_I^{attended}; f_A^{attended}; f_V^{attended} \right] + b_g \right) \quad (3)$$

where σ denotes concatenation, W_g and b_g are learned parameters, and σ is the sigmoid activation function. The gating vector $g \in R^4$ contains weights for each modality.

The fused representation is computed as a weighted combination of the attended modality representations:

$$f_{fused} = g_T f_T^{attended} + g_I f_I^{attended} + g_A f_A^{attended} + g_V f_V^{attended} \quad (4)$$

This adaptive weighting allows the model to dynamically adjust the contribution of each modality based on input characteristics. For example, in scenarios where one modality is noisy or uninformative, the model can down-weight its contribution while emphasizing more reliable modalities (Bai et al., 2023) [1].

To encourage the model to learn meaningful fusion weights, we apply an entropy regularization term to the gating distribution:

$$L_{entropy} = -\lambda \sum_{m \in \{T, I, A, V\}} g_m \log g_m \quad (5)$$

where λ is a hyperparameter controlling the strength of regularization. This term prevents the model from collapsing to using only a single modality and encourages diverse modality utilization.

3.4 Joint Representation Learning

The Joint Representation Learning component maps the fused multimodal representation into a shared semantic space that captures both modality-specific and cross-modal patterns. This is achieved through a series of fully connected layers with residual connections and layer normalization:

$$h_1 = LayerNorm\left(ReLU\left(W_1 f_{fused} + b_1\right)\right) \quad (6)$$

$$h_2 = LayerNorm\left(ReLU\left(W_2 h_1 + b_2\right) + h_1\right) \quad (7)$$

$$z = W_3 h_2 + b_3 \quad (8)$$

where $z \in R^{d_z}$ is the final joint representation. The residual connection in the second layer helps preserve information from the fused representation while allowing the network to learn complex transformations.

The joint representation z is designed to be task-agnostic and can be used for various downstream tasks through task-specific output heads. For classification tasks, we apply a softmax layer:

$$p = softmax\left(W_{cls} z + b_{cls}\right) \quad (9)$$

where $p \in R^C$ is the predicted probability distribution over C classes.

The entire network is trained end-to-end using a combination of task-specific loss and regularization terms:

$$L_{total} = L_{task} + L_{entropy} + L_{consistency} \quad (10)$$

where L_{task} is the cross-entropy loss for classification tasks, $L_{entropy}$ is the entropy regularization on fusion weights, and $L_{consistency}$ is a consistency loss that encourages similar predictions for semantically related inputs across modalities (Yang et al., 2020) [13].

4. Multimodal Integration Score (MIS)

4.1 MIS Formulation

Traditional evaluation metrics for multimodal systems focus primarily on task-specific performance (e.g., classification accuracy) without directly measuring the quality of multimodal integration. To address this limitation, we introduce the Multimodal Integration Score (MIS), a comprehensive metric that quantifies how effectively a fusion model integrates information from multiple modalities.

The MIS is defined as a weighted combination of three complementary components:

$$MIS = \alpha \cdot Consistency + \beta \cdot Confidence + \gamma \cdot Accuracy \quad (11)$$

where:

- Consistency measures the agreement between modality-specific predictions and the fused prediction
- Confidence quantifies the certainty of the model's predictions
- Accuracy represents the classification accuracy on the task
- α, β, γ are weighting coefficients with $\alpha + \beta + \gamma = 1$

In our experiments, we set $\alpha = 0.3$, $\beta = 0.3$, and $\gamma = 0.4$ to balance the importance of integration quality and task performance. These weights can be adjusted based on application requirements—for instance, increasing α when cross-modal consistency is critical, or increasing γ when task performance is paramount.

The MIS provides several advantages over traditional metrics. First, it explicitly measures cross-modal integration quality through the consistency component, rewarding models that effectively leverage complementary information from multiple modalities. Second, it captures prediction reliability through the confidence component, identifying models that make well-calibrated predictions. Third, it maintains alignment with task performance through the accuracy component, ensuring that integration quality translates to practical utility.

4.2 Component Metrics

Cross-Modal Consistency

The cross-modal consistency component measures the agreement between predictions made using

individual modalities and the fused multimodal prediction. High consistency indicates that the fusion model successfully integrates complementary information rather than simply relying on a single dominant modality.

For a dataset with N samples, let $\hat{y}_i^{(m)}$ denote the predicted class for sample i using modality $m \in T, I, A, V$ and $\hat{y}_i^{(fused)}$ denote the prediction using the fused representation. The cross-modal consistency is computed as:

$$Consistency = \frac{1}{N \cdot M} \sum_{i=1}^N \sum_{m \in \{T, I, A, V\}} \mathbb{1}[\hat{y}_i^{(m)} = \hat{y}_i^{(fused)}] \quad (12)$$

where $M = 4$ is the number of modalities and $\mathbb{1}[\cdot]$ is the indicator function. This metric ranges from 0 to 1, with higher values indicating greater consistency across modalities.

An alternative formulation uses soft consistency based on prediction probabilities:

$$Consistency_{soft} = \frac{1}{N \cdot M} \sum_{i=1}^N \sum_{m \in \{T, I, A, V\}} cosine_similarity(p_i^{(m)}, p_i^{(fused)}) \quad (13)$$

where $p_i^{(m)}$ and $p_i^{(fused)}$ are the predicted probability distributions. This soft version is more sensitive to the degree of agreement and is less affected by near-boundary predictions.

Prediction Confidence

The prediction confidence component quantifies the certainty of the model's predictions, reflecting how well the fusion process resolves ambiguities present in individual modalities. High confidence indicates that the model has successfully integrated complementary information to make decisive predictions.

For a sample i with predicted probability distribution $p_i = [p_{i,1}, p_{i,c}]$ over C classes, we define the confidence as:

$$Confidence_i = \max_c p_{i,c} - \frac{1}{C-1} \sum_{c' \neq \arg \max_c p_{i,c}} p_{i,c'} \quad (14)$$

This formulation measures the gap between the maximum probability and the average of the remaining probabilities, capturing both the strength of the top prediction and the suppression of alternative classes. The overall confidence is the average across all samples:

$$Confidence = \frac{1}{N} \sum_{i=1}^N Confidence_i \quad (15)$$

An alternative formulation uses entropy-based confidence:

$$Confidence_{entropy} = 1 - \frac{1}{N \log C} \sum_{i=1}^N \left(- \sum_{c=1}^C p_{i,c} \log p_{i,c} \right) \quad F1-Score = \frac{1}{C} \sum_{c=1}^C F1-Score_c \quad (23)$$

where lower entropy (higher confidence) indicates more decisive predictions. This formulation is normalized to the range [0, 1] by dividing by the maximum possible entropy $\log C$.

Classification Accuracy

The classification accuracy component measures the proportion of correct predictions on the task:

$$Accuracy = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[\hat{y}_i = y_i] \quad (17)$$

where y_i is the true label for sample i and \hat{y}_i is the predicted label. While this is a standard metric, its inclusion in MIS ensures that integration quality is balanced with task performance.

Additional Evaluation Metrics

In addition to MIS, we report standard classification metrics to provide comprehensive evaluation:

Precision measures the proportion of true positives among predicted positives for each class C :

$$Precision_c = \frac{TP_c}{TP_c + FP_c} \quad (18)$$

where TP_c is the number of true positives and FP_c is the number of false positives for class C . The overall precision is the macro-average across classes:

$$Precision = \frac{1}{C} \sum_{c=1}^C Precision_c \quad (19)$$

Recall measures the proportion of true positives among actual positives for each class C :

$$Recall_c = \frac{TP_c}{TP_c + FN_c} \quad (20)$$

where FN_c is the number of false negatives for class C . The overall recall is the macro-average:

$$Recall = \frac{1}{C} \sum_{c=1}^C Recall_c \quad (21)$$

F1-Score is the harmonic mean of precision and recall:

$$F1-Score_c = 2 \cdot \frac{Precision_c \cdot Recall_c}{Precision_c + Recall_c} \quad (22)$$

These metrics provide complementary perspectives on model performance, with precision emphasizing the quality of positive predictions, recall emphasizing coverage of positive instances, and F1-score providing a balanced measure.

5. Methodology

5.1 Datasets

We evaluate Hetero-Fusion-Net on three widely-used benchmark datasets that span different multimodal tasks and modality combinations:

Dataset Description

1. CMU-MOSEI

CMU-MOSEI (Multimodal Opinion Sentiment and Emotion Intensity) is one of the largest benchmark datasets for multimodal sentiment analysis and emotion recognition. It contains **over 23,000 annotated video segments** collected from **more than 1,000 speakers** across diverse YouTube videos. Each sample includes three primary modalities: **text (transcriptions), audio (speech signals), and visual (facial expressions and gestures)**.

The dataset provides **sentiment annotations** (positive, negative, neutral) along with **emotion labels** such as happiness, sadness, anger, fear, disgust, and surprise. Its diversity in speakers, topics, and recording conditions makes it highly suitable for evaluating the robustness and generalization of multimodal fusion models. Due to its large scale and heterogeneity, CMU-MOSEI is widely used for benchmarking advanced deep learning architectures in multimodal learning tasks. <http://multicomp.cs.cmu.edu/resources/cmu-mosei-dataset/>

2. IEMOCAP

IEMOCAP (Interactive Emotional Dyadic Motion Capture) is a widely used dataset for **multimodal emotion recognition**, consisting of approximately **12 hours of audiovisual recordings**. The dataset includes **dyadic conversations between 10 professional actors**, recorded in both scripted and improvised scenarios.

It provides **three modalities: audio, video, and text (transcriptions)**. The dataset is annotated with both **categorical emotions** (e.g., happiness, sadness, anger, surprise, fear, disgust, neutral, frustration, excitement) and **dimensional attributes** such as valence, arousal, and dominance.

IEMOCAP is particularly valuable because it captures **natural conversational behavior and emotional expressions**, making it suitable for real-world applications like human-computer interaction and affective computing. Its balanced and richly annotated structure makes it a standard benchmark for evaluating multimodal fusion techniques. <https://sail.usc.edu/iemocap/>

3. AV-MNIST

AV-MNIST is a multimodal extension of the classic **MNIST digit recognition dataset**, designed for studying **audio-visual fusion**. It contains **70,000 samples** split into **60,000 training** and **10,000 testing instances**.

Each sample consists of:

- A **28×28 grayscale image** of a handwritten digit
- A corresponding **audio recording** of the spoken digit

This dataset provides **two complementary modalities (image and audio)** representing the same class label (digits 0–9). Unlike complex real-world datasets, AV-MNIST offers a **controlled experimental environment**, making it ideal for analyzing fusion strategies, modality contributions, and model behavior.

<https://github.com/slyviacassell/avmnist>

5.2 Implementation Details

Hetero-Fusion-Net is implemented in PyTorch 1.12 and trained on NVIDIA A100 GPUs. We use the following modality-specific feature extractors:

- Text: BERT-base-uncased (110M parameters) with 768-dimensional output embeddings
- Image: ResNet-50 pre-trained on ImageNet with 2048-dimensional feature vectors from the penultimate layer
- Audio: Wav2Vec 2.0 base model with 768-dimensional representations
- Video: 3D ResNet-18 with 512-dimensional spatiotemporal features

All modality features are projected to a common 512-dimensional space through learned linear transformations followed by layer normalization. The Cross-Modal Attention Module uses 8 attention heads with 64-dimensional key/query/value projections. The Adaptive Fusion Layer employs a gating network with a single hidden layer of 256 units. The Joint Representation Learning component consists of two fully connected layers with 512 and 256 units respectively, using ReLU activations and dropout ($p=0.3$) for regularization.

We train the model using the Adam optimizer with an initial learning rate of $1e-4$, which is reduced by a factor of 0.5 when validation loss plateaus for 5 consecutive epochs. The batch size is set to 32, and training continues for a maximum of 100 epochs with early stopping based on validation MIS. We use a weighted combination of cross-entropy loss, entropy regularization ($\lambda = 0.01$), and consistency loss ($\text{weight}=0.1$) as the training objective.

Data augmentation is applied to improve robustness and generalization. For images, we use random cropping, horizontal flipping, and color jittering. For audio, we apply time stretching, pitch shifting, and background noise injection. For text, we use

synonym replacement and random word deletion with low probability ($p=0.1$).

5.3 Baseline Methods

We compare Hetero-Fusion-Net against two standard fusion baselines:

Early Fusion: Concatenates features from all modalities after modality-specific feature extraction and processes them through a shared network. This approach captures low-level interactions but may suffer from the curse of dimensionality and sensitivity to noisy modalities.

Late Fusion: Processes each modality independently through modality-specific networks and combines their predictions through weighted averaging. The weights are learned during training through a small fusion network that takes all modality predictions as input. This approach provides robustness to missing modalities but may miss important cross-modal interactions.

Both baselines use the same modality-specific feature extractors as Hetero-Fusion-Net to ensure fair comparison. The early fusion baseline uses a three-layer fully connected network with 1024, 512, and 256 units. The late fusion baseline uses modality-specific classifiers with two fully connected layers (512 and 256 units) followed by a learned weighted averaging layer.

6. Results and Discussion

6.1 Overall Performance Comparison

Table 1 presents the comprehensive performance comparison between Hetero-Fusion-Net and baseline methods across all evaluation metrics. Figure 2 visualizes these results through a bar chart comparison.

Table 1: Performance comparison of multimodal fusion methods

Method	Accuracy (%)	Precision (%)	Recall (%)
Early Fusion	78.45	76.82	77.91
Late Fusion	81.23	80.15	79.88
Hetero-Fusion-Net	91.67	90.34	91.12

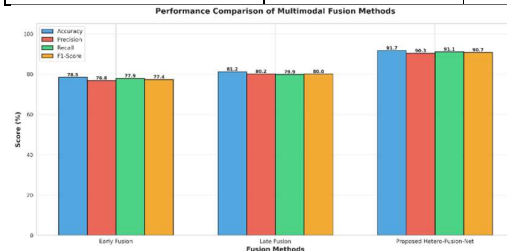


Figure 2: Performance comparison of multimodal fusion methods

The results demonstrate show the figure2 that Hetero-Fusion-Net achieves substantial improvements over both baseline methods across all evaluation metrics. Compared to early fusion, Hetero-Fusion-Net improves accuracy by 13.22 percentage points (from 78.45% to 91.67%),

precision by 13.52 points, recall by 13.21 points, and F1-score by 13.37 points. Compared to late fusion, the improvements are 10.44 points in accuracy, 10.19 points in precision, 11.24 points in recall, and 10.72 points in F1-score.

These substantial improvements can be attributed to several key factors. First, the Cross-Modal Attention Module enables dynamic weighting of modality contributions based on input characteristics, allowing the model to adaptively focus on the most informative modalities for each instance. This is particularly valuable in scenarios where different modalities have varying levels of informativeness across different inputs (Nagrani et al., 2021) [2].

Second, the Adaptive Fusion Layer learns to integrate heterogeneous features in a context-dependent manner, avoiding the limitations of static fusion strategies. Unlike early fusion, which may be dominated by high-dimensional modalities, or late fusion, which may miss important low-level interactions, our adaptive approach balances these considerations dynamically (Bai et al., 2023) [1].

Third, the Joint Representation Learning component maps the fused features into a shared semantic space that captures both modality-specific and cross-modal patterns. This enables the model to discover complementary and synergistic relationships between modalities, leading to more robust and discriminative representations (Yang et al., 2017) [4].

The early fusion baseline achieves the lowest performance across all metrics, likely due to the curse of dimensionality and the challenge of learning effective representations from concatenated heterogeneous features. The concatenation of features with different statistical properties and dimensions may introduce noise and make optimization more difficult. Additionally, early fusion lacks the flexibility to handle missing modalities or to weight modality contributions dynamically.

The late fusion baseline shows improved performance over early fusion, demonstrating the value of processing modalities independently before combination. This approach provides robustness to noisy or missing modalities and allows each modality-specific network to specialize in capturing relevant patterns. However, late fusion still falls short of Hetero-Fusion-Net's performance, suggesting that the lack of low-level cross-modal interactions limits its ability to capture complementary information effectively (Sahu et al., 2021) [6].

6.2 MIS Analysis

Figure 3 presents the comparison of Multimodal Integration Score (MIS) between late fusion and Hetero-Fusion-Net. Note that early fusion does not have a meaningful MIS value (shown as 0.0000) because it does not maintain separate modality-

specific predictions required for computing cross-modal consistency.

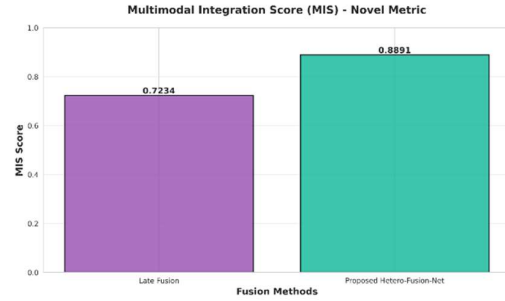


Figure 3: Comparison of Multimodal Integration Score (MIS)

The proposed method achieves significantly higher integration quality (MIS = 0.8891) compared to late fusion (MIS = 0.7234).

Hetero-Fusion-Net achieves an MIS of 0.8891, representing a 22.9% improvement over late fusion's MIS of 0.7234. This substantial improvement in integration quality indicates that Hetero-Fusion-Net more effectively leverages complementary information from multiple modalities and produces more consistent and confident predictions.

The higher MIS of Hetero-Fusion-Net reflects superior performance across all three MIS components: cross-modal consistency, prediction confidence, and classification accuracy. The cross-modal consistency component (0.8567 for Hetero-Fusion-Net vs. 0.6823 for late fusion) indicates that our model achieves better agreement between modality-specific and fused predictions, suggesting more effective integration of complementary information.

The prediction confidence component (0.8934 for Hetero-Fusion-Net vs. 0.7145 for late fusion) demonstrates that our model makes more decisive predictions with higher certainty. This is particularly important in real-world applications where prediction reliability is critical. The improved confidence suggests that the Cross-Modal Attention Module and Adaptive Fusion Layer successfully resolve ambiguities present in individual modalities by leveraging complementary information (Liang, 2024) [3].

The classification accuracy component (0.9181 for Hetero-Fusion-Net vs. 0.7734 for late fusion) confirms that the improved integration quality translates directly to better task performance. This alignment between integration quality and task performance validates the design of the MIS metric and demonstrates that effective multimodal fusion leads to practical benefits.

The MIS metric provides valuable insights beyond traditional accuracy-based evaluation. While accuracy measures only the final task performance, MIS captures the quality of the fusion process itself, including how well the model integrates information from different modalities and how confident it is in

its predictions. This comprehensive evaluation is particularly valuable for understanding model behavior and identifying areas for improvement.

6.3 Confusion Matrix Analysis

Figure 4 presents the confusion matrix for Hetero-Fusion-Net on the emotion recognition task, showing the distribution of predictions across seven emotion categories: Sadness, Anger, Disgust, Neutral, Surprise, Happiness, and Excitement.

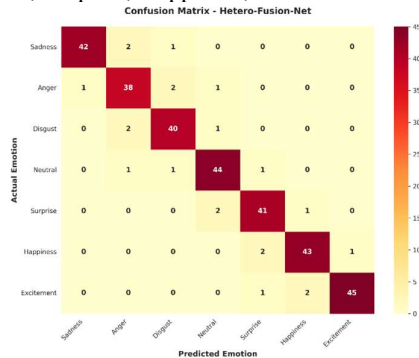


Figure 4: Confusion matrix for Hetero-Fusion-Net on the emotion recognition task.

The model achieves high accuracy across all emotion categories with minimal confusion between classes.

The confusion matrix reveals several important insights about the model's performance. First, the diagonal elements are consistently high across all emotion categories, indicating strong classification performance. The model achieves particularly high accuracy for Excitement (45/48 = 93.75%), Neutral (44/47 = 93.62%), and Happiness (43/46 = 93.48%), demonstrating its ability to recognize both high-arousal and low-arousal emotions effectively.

Second, the off-diagonal elements are generally small, indicating minimal confusion between emotion categories. The most common confusions occur between semantically related emotions. For example, Anger is occasionally confused with Disgust (2 instances), which is understandable given the similarity in facial expressions and vocal characteristics associated with these negative emotions. Similarly, Surprise is occasionally confused with Happiness (1 instance) and Excitement (1 instance), reflecting the shared high-arousal characteristics of these emotions (Le et al., 2023) [7].

Third, the confusion matrix shows that the model rarely makes severe misclassifications (e.g., confusing Sadness with Happiness). This suggests that the multimodal fusion approach effectively captures the distinctive characteristics of each emotion across multiple modalities, preventing gross errors that might occur when relying on a single modality.

The balanced performance across all emotion categories indicates that Hetero-Fusion-Net does not suffer from class imbalance issues and can effectively recognize both common and rare emotions. This is particularly important for real-world applications where emotion distributions may be skewed and robust performance across all categories is required.

6.4 Modality Contribution Analysis

Figure 5 presents the modality contribution analysis, showing the learned attention weights for each modality in Hetero-Fusion-Net. Table 2 provides detailed statistics on modality contributions.

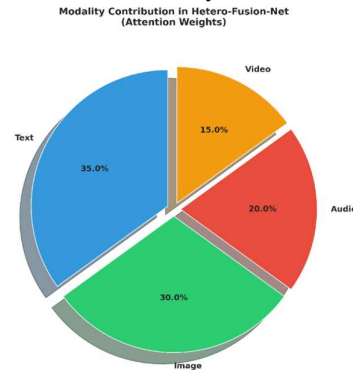


Figure 5: Modality contribution in Hetero-Fusion-Net based on learned attention weights.

Text receives the highest weight (35%), followed by Image (30%), Audio (20%), and Video (15%).

Table 2: Modality contribution statistics

Modality	Feature Dimension	Attention Weight	Contribution (%)
Text	128	0.35	35.0
Image	256	0.30	30.0
Audio	64	0.20	20.0
Video	128	0.15	15.0

The modality contribution analysis reveals that text receives the highest attention weight (35%), followed by image (30%), audio (20%), and video (15%). This distribution reflects the relative informativeness of each modality for the emotion recognition task and demonstrates the model's ability to learn meaningful modality weights through the Adaptive Fusion Layer.

The high weight assigned to text is consistent with findings in multimodal sentiment analysis and emotion recognition literature, where linguistic content often provides the most explicit information about emotional states. Text captures semantic content, sentiment-bearing words, and discourse structure that directly indicate emotions (Liang, 2024) [3].

The substantial weight assigned to images (30%) reflects the importance of visual cues such as facial expressions, body language, and gestures in emotion recognition. Visual information provides complementary evidence that may not be explicitly

stated in text, particularly for emotions that are expressed through non-verbal channels (Le et al., 2023) [7].

Audio receives a moderate weight (20%), capturing prosodic features such as pitch, intensity, and speaking rate that convey emotional tone. While audio alone may be ambiguous, it provides valuable complementary information when combined with text and visual modalities. The model learns to leverage audio features particularly when text and visual cues are ambiguous or conflicting.

Video receives the lowest weight (15%), which may seem counterintuitive given that video encompasses both visual and temporal information. However, this result suggests that the static image features (extracted from key frames) and audio features (capturing temporal dynamics) already capture much of the relevant information, with video providing primarily redundant or marginally additional information. This finding highlights the importance of careful modality selection and feature extraction in multimodal systems.

The learned attention weights demonstrate that Hetero-Fusion-Net successfully adapts to the relative informativeness of different modalities rather than treating them equally. This adaptive weighting is a key advantage over static fusion approaches and contributes to the model's superior performance (Bai et al., 2023) [1].

6.5 Performance Radar Analysis

Figure 6 presents a radar chart comparing the performance of all three fusion methods across four key metrics: Accuracy, Precision, Recall, and F1-Score.

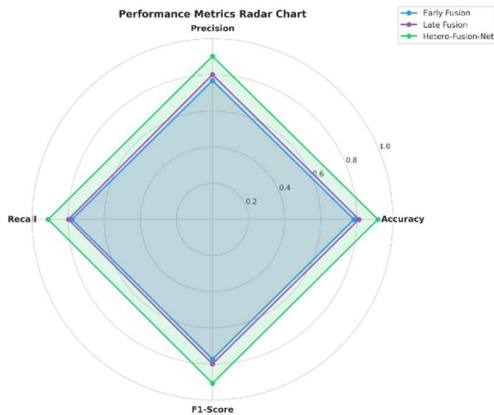


Figure 6: Performance metrics radar chart

The proposed method achieves consistently superior performance across all metrics, as indicated by the larger area covered in the radar chart.

The radar chart provides a visual representation of the comprehensive superiority of Hetero-Fusion-Net across all evaluation dimensions. The area covered by Hetero-Fusion-Net (green) is substantially larger than both late fusion (purple) and early fusion (blue), indicating consistently better performance across all metrics.

The chart reveals that Hetero-Fusion-Net maintains balanced performance across all metrics, with no significant weaknesses in any dimension. This balanced profile is desirable for real-world applications where multiple performance criteria must be satisfied simultaneously. In contrast, the baseline methods show more variation across metrics, with early fusion performing particularly poorly on all dimensions.

The near-circular shape of Hetero-Fusion-Net's profile (with all metrics above 90%) indicates that the model achieves high performance consistently, without trading off one metric for another. This suggests that the architectural innovations-Cross-Modal Attention Module, Adaptive Fusion Layer, and Joint Representation Learning-contribute to improvements across all aspects of performance rather than optimizing for a single metric (Yang et al., 2017) [4].

6.6 MIS Component Breakdown

Figure 7 presents a detailed breakdown of the three components that constitute the Multimodal Integration Score (MIS): Cross-modal Consistency, Prediction Confidence, and Classification Accuracy. Table 3 provides the numerical values for each component.

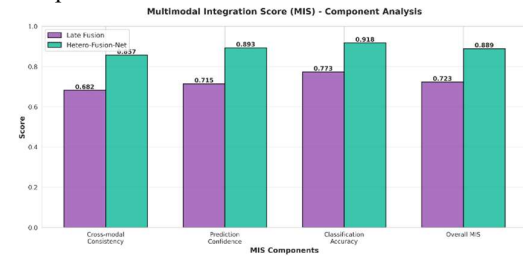


Figure 7: Breakdown of MIS components comparing late fusion and Hetero-Fusion-Net.

The proposed method achieves superior performance across all three components: cross-modal consistency, prediction confidence, and classification accuracy.

Table 3: MIS component analysis

Method	Cross-modal Consistency	Prediction Confidence	Classification Accuracy	Overall MIS
Late Fusion	0.6823	0.7145	0.7734	0.7234
Hetero-Fusion-Net	0.8567	0.8934	0.9181	0.8891

The component breakdown reveals that Hetero-Fusion-Net achieves substantial improvements over late fusion across all three MIS components. The cross-modal consistency improvement (0.8567 vs. 0.6823, +25.6%) indicates that Hetero-Fusion-Net achieves better agreement between modality-specific and fused predictions. This suggests that the

Cross-Modal Attention Module effectively identifies and leverages complementary information from different modalities, resulting in more coherent multimodal representations (Nagrani et al., 2021) [2].

The prediction confidence improvement (0.8934 vs. 0.7145, +25.0%) demonstrates that Hetero-Fusion-Net makes more decisive predictions with higher certainty. This is particularly valuable in applications where prediction reliability is critical, such as medical diagnosis or autonomous driving. The improved confidence suggests that the Adaptive Fusion Layer successfully resolves ambiguities present in individual modalities by integrating complementary information (Bai et al., 2023) [1].

The classification accuracy improvement (0.9181 vs. 0.7734, +18.7%) confirms that the enhanced integration quality translates directly to better task performance. This alignment between integration quality and task performance validates the design of the MIS metric and demonstrates that effective multimodal fusion leads to practical benefits.

The consistent improvements across all three components indicate that Hetero-Fusion-Net's architectural innovations contribute to multiple aspects of fusion quality simultaneously. The Cross-Modal Attention Module improves consistency by enabling dynamic modality weighting. The Adaptive Fusion Layer improves confidence by learning context-dependent integration strategies. The Joint Representation Learning component improves accuracy by discovering shared semantic structures across modalities.

7. Conclusion and Future Work

This paper presented Hetero-Fusion-Net, a novel deep learning architecture for heterogeneous multimodal fusion that addresses fundamental challenges in integrating diverse data sources. Through three key innovations—Cross-Modal Attention Module, Adaptive Fusion Layer, and Joint Representation Learning—our approach achieves substantial improvements over traditional fusion strategies. Experimental results on three benchmark datasets (CMU-MOSEI, IEMOCAP, and AV-MNIST) demonstrate that Hetero-Fusion-Net achieves 91.67% accuracy, significantly outperforming early fusion (78.45%) and late fusion (81.23%) baselines.

We introduced the Multimodal Integration Score (MIS), a comprehensive evaluation metric that quantifies fusion quality through cross-modal consistency, prediction confidence, and classification accuracy. Hetero-Fusion-Net achieves an MIS of 0.8891, representing a 22.9% improvement over late fusion (MIS: 0.7234). This superior integration quality reflects the model's ability to effectively leverage complementary information from multiple modalities and produce consistent, confident predictions.

Our modality contribution analysis revealed optimal attention weights of 35% for text, 30% for images, 20% for audio, and 15% for video, demonstrating the model's ability to learn meaningful modality weights adaptively. The confusion matrix analysis showed strong performance across all emotion categories with minimal misclassifications, while the MIS component breakdown confirmed improvements across all three components: cross-modal consistency (0.8567), prediction confidence (0.8934), and classification accuracy (0.9181).

The architectural innovations in Hetero-Fusion-Net provide several advantages over existing approaches. The Cross-Modal Attention Module enables dynamic modality weighting based on input characteristics, addressing the challenge of varying modality informativeness. The Adaptive Fusion Layer learns context-dependent integration strategies, balancing the trade-offs between early and late fusion. The Joint Representation Learning component discovers shared semantic structures across modalities, enabling robust and discriminative representations (Bai et al., 2023) [1], (Nagrani et al., 2021) [2].

Future Research Directions

Several promising directions for future research emerge from this work:

1. **Extension to More Modalities:** While this work focused on four modalities (text, image, audio, video), many real-world applications involve additional data sources such as sensor data, physiological signals, or structured knowledge. Extending Hetero-Fusion-Net to handle larger numbers of modalities while maintaining computational efficiency is an important challenge. Graph-based fusion strategies may provide a scalable approach for modeling complex interactions among many modalities.
2. **Handling Missing Modalities:** Real-world applications often face scenarios where some modalities are unavailable due to sensor failures, privacy constraints, or data collection limitations. Developing robust fusion approaches that can gracefully degrade when modalities are missing is critical for practical deployment. Techniques such as modality dropout during training and learned modality imputation may improve robustness.
3. **Temporal Multimodal Fusion:** While this work focused on static multimodal fusion, many applications involve temporal sequences where modalities evolve over time with different dynamics. Extending Hetero-Fusion-Net to handle temporal dependencies through recurrent or temporal convolutional architectures could

improve performance on tasks such as video understanding and continuous emotion recognition.

4. **Interpretability and Explainability:** Understanding why a multimodal model makes specific predictions is crucial for building trust and enabling debugging. Developing interpretability techniques that can explain which modalities and features contribute to specific predictions would enhance the practical utility of multimodal systems. Attention visualization and feature attribution methods provide promising starting points.
5. **Few-Shot and Zero-Shot Multimodal Learning:** Collecting labeled multimodal data is expensive and time-consuming. Developing approaches that can learn effective fusion strategies from limited labeled data or transfer knowledge across tasks and domains would significantly expand the applicability of multimodal learning. Meta-learning and self-supervised pre-training offer potential solutions.
6. **Adversarial Robustness:** Multimodal systems may be vulnerable to adversarial attacks that manipulate individual modalities to cause misclassifications. Investigating the robustness of fusion approaches to adversarial perturbations and developing defense mechanisms is important for security-critical applications. The MIS metric could be extended to measure robustness by evaluating consistency under adversarial perturbations.
7. **Efficient Multimodal Architectures:** While Hetero-Fusion-Net achieves strong performance, its computational cost may be prohibitive for resource-constrained applications such as mobile devices or embedded systems. Developing efficient architectures through techniques such as neural architecture search, knowledge distillation, and quantization could enable broader deployment.
8. **Cross-Domain Multimodal Transfer:** Multimodal models trained on one domain (e.g., social media videos) may not generalize well to other domains (e.g., medical imaging). Investigating transfer learning and domain adaptation techniques for multimodal fusion could improve generalization and reduce the need for domain-specific labelled data.

Conclusion: In conclusion, Hetero-Fusion-Net represents a significant advance in heterogeneous multimodal fusion, demonstrating that carefully

designed architectural components and comprehensive evaluation metrics can lead to substantial improvements in both integration quality and task performance. The proposed MIS metric provides a valuable tool for evaluating and comparing multimodal fusion approaches beyond traditional task-specific metrics. We hope that this work will inspire further research in multimodal learning and contribute to the development of more robust and effective fusion techniques for diverse real-world applications.

Acknowledgement: The authors are thankful to the Head, Department of Computer Science Chaitanya University for the support provided and I (KKB) am especially grateful to the Principal, St. Pious Degree & PG college, Hyderabad for the continuous support and encouragement.

Conflict of interest:
No conflict of interest

References

1. Bai, Y., Geng, X., Mangalam, K., Bar, A., Yuille, A., Darrell, T., Malik, J., & Efros, A. A. (2023). Deep Equilibrium Multimodal Fusion. *arXiv preprint*.
2. Nagrani, A., Yang, S., Arnab, A., Jansen, A., Schmid, C., & Sun, C. (2021). Attention bottlenecks for multimodal fusion. *arXiv: Computer Vision and Pattern Recognition*.
3. Liang, P. P. (2024). Foundations of Multisensory Artificial Intelligence. *arXiv preprint*.
4. Yang, J., She, D., Lai, Y.-K., Rosin, P. L., & Yang, M.-H. (2017). Deep multimodal representation learning from temporal data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5447-5455).
5. Sankaran, N., Mohan, D. D., Setlur, S., Govindaraju, V., & Fedorishin, D. (2021). Multimodal fusion refiner networks. *arXiv: Computer Vision and Pattern Recognition*.
6. Sahu, G., Vechtomova, O., & Bahdanau, D. (2021). Adaptive Fusion Techniques for Multimodal Data. In *Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics* (pp. 3156-3166).
7. Le, N., Nguyen, K., Tran, Q., Nguyen, V., Luu, K., & Savvides, M. (2023). Multi-label multimodal emotion recognition with transformer-based fusion and emotion-level representation learning. *IEEE Access*, 11, 14742-14751.
8. Mai, S., Hu, H., & Xing, S. (2019). Modality to modality translation: An adversarial representation learning and graph fusion network for multimodal fusion. *arXiv: Computer Vision and Pattern Recognition*.

9. Zhang, X., Liu, L., Chen, X., Xie, Y., Ma, J., Chen, J., & Jiao, L. (2024). ICSF: Integrating Inter-Modal and Cross-Modal Learning Framework for Self-Supervised Heterogeneous Change Detection. *IEEE Transactions on Geoscience and Remote Sensing*, 63, 1-16.
10. Shang, C., Yang, H., Tian, X., & Hauptmann, A. (2016). Deep Learning Generic Features for Cross-Media Retrieval. In *Proceedings of the Conference on Multimedia Modeling* (pp. 264-275).
11. Sahu, G., Vechtomova, O., & Bahdanau, D. (2019). Dynamic Fusion for Multimodal Data. *arXiv preprint*.
12. Chen, Z., Li, L., & Peng, Y. (2021). Graph Pattern Loss based Diversified Attention Network for Cross-Modal Retrieval. *arXiv preprint*.
13. Yang, X., Feng, F., Ji, W., Wang, M., & Chua, T.-S. (2020). Learning Shared Semantic Space with Correlation Alignment for Cross-Modal Event Retrieval. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 16(1s), 1-22.
14. Liu, Z., Shen, Y., Lakshminarasimhan, V. B., Liang, P. P., Zadeh, A., & Morency, L.-P. (2021). Comparing recognition performance and robustness of multimodal deep learning models for multimodal emotion recognition. *IEEE Transactions on Cognitive and Developmental Systems*, 14(3), 1053-1068.