# Efficient Deep Q Learning Models Featuring Adaptive Attention and Emotionally Augmented Fusion for Multimodal Sentiment Analysis

Anupam Chaube; Dr. Usha  Kosarkar
Department  of  Science  &  Technology
G.H.Raisoni College of Engineering & Management Nagpur,India


Devarshi A. Patrikar
Department  of  Science  &  Technology
G.H.Raisoni  Skill  Tech  University  Nagpur,India

**Abstract**— The interplay of the features of text, audio, and video makes multimodal sentiment analysis very critical in psychological assessments. Traditional methods have failed to prioritize relevant features well in order to capture the rich emotional dynamics that are central to most psychological insights. Most limitations associated with existing methods include not having a proper modality-specific focus on relevance and poor robustness to noisy inputs along with the lack of interpretability sets. In this direction, we introduce an integrated framework with three novel deep Q-learning models that could enhance multimodal sentiment analysis. These are Adaptive Attention-Gated Deep Q-Learning (AADQL), Emotionally Augmented Reinforcement Fusion Network (EARFN), and Hierarchical Task-Aware Deep Q- Learning (HTADQL). AADQL incorporates the adaptive attention gating mechanism in order to dynamically weigh the relevance of features across the different modalities, and we found it to result in improved accuracy of 4-6% and reducing computational overhead by 15-20%. EARFN employs psychological emotion models to combine emotion- specific representations within a reinforcement learning paradigm, and it produces high robustness with 25-30% noise tolerance and superior prediction accuracy of 88- 90%. HTADQL makes use of a hierarchical architecture, which solves the subtasks of valence and arousal before combining them for sentiment classification, thus showing better interpretability and error rate reduction by 30-40%. These methodologies hold some significant implications: they improve the accuracy up to 10% against baseline models; exhibit robustness in noisy environments and deliver increased interpretability about the psychological sentiment dynamics. The work that is reported here brings in a remarkable advancement on the deep Q-learning front for multimodal sentiment analysis and would be even more apt to create room for even more efficacious psychological assessment tools.

**Keywords—**
Deep Q Learning, Multimodal Sentiment Analysis, Adaptive Attention, Emotionally Augmented Fusion, Psychological Assessments

## I.Introduction

For psychological assessments, accurate sentiment analysis is imperative because insights drawn from multimodal data, text, audio, and video, make a comprehensive understanding of human emotions and behaviors possible. The traditional methods usually suffer from suboptimal feature prioritization, weak robustness against noisy inputs, and low interpretability sets. Such drawbacks prevent these methods from being used where the need is for the subtlety and accuracy of psychological insights in a process. AADQL introduces an adaptive attention gating mechanism for dynamic weighing of relevant features across modalities with significantly improved accuracy by 4-6% and reduction in computational overhead of 15-20%. EARFN uses psychological emotion models for the combination of emotion-specific representations within a reinforcement learning paradigm that results in high robustness with 25- 30% noise tolerance and superior prediction accuracy of 88- 90%. HTADQL has a hierarchical architecture, solving subtasks like valence and arousal before combining them for sentiment classification that results in better interpretability and 30-40% error rate reduction. The methods are impactful: they boost accuracy up to 10% relative to baseline models, and their robustness in noisy environments and increased interpretability of psychological sentiment dynamics can be achieved.

The framework represented here marks a significant advancement in deep Q-learning for multimodal sentiment analysis that would pave the way to more effective psychological assessment tools. In the case of psychological assessments, the sentiment analysis has to be precise since the insights gathered from multimodal data - text, audio, and video - enable an overall understanding of human emotions and behaviors.

The traditional methods usually suffer from suboptimal feature prioritization, weak robustness against noisy inputs, and low interpretability sets. Such drawbacks prevent these methods from being used where the need is for the subtlety and accuracy of psychological insights in a process. To overcome these issues [1, 2, 3], this work provides an advanced deep reinforcement learning framework using novel architectures in multimodal sentiment analysis. The Adaptive Attention-Gated Deep Q-Learning (AADQL) model assigns modality-specific weights dynamically through adaptive attention gates, which enables feature prioritization with precision and efficiency. Additionally, the Emotionally Augmented Reinforcement Fusion Network (EARFN) integrates emotion-specific representations inspired by psychological emotion models, significantly improving noise robustness and accuracy. Last but not least, the Hierarchical Task-Aware Deep Q-Learning framework adopts the hierarchical structure to handle all subtasks like valence and arousal prediction before achieving the final classification with stronger interpretability and fine-grained learning. Therefore, this work advances a cutting-edge technique for robust and fine-tuning precision in the accuracy of predicting sentiment by utilizing methodologies that fill up some gaps with the ones presently used. This proposal has new standards of applying deep Q-learning to multimodal sentiment analysis, opening the doorway for reliable and interpretable tools in psychological assessment. They could achieve up to 10% accuracy improvement and showed robust performance under noisy conditions. There is considerable recent interest in the integration of machine learning with multimodal sentiment analysis applied to mental health and psychological assessments process.

## A. Review of Existing Models Sentiment Analysis

The work presents key studies, a detailed survey on state-of- the-art approaches with their implications for the process. Lopes et al. [1] developed a novel dataset called PerceptSent, where convolutional neural networks (CNNs) are used and coupled with the external knowledge application of visual sentiment analysis. The research focused on deep networks that allow for the probing of subjectivity in multimodal data, thereby providing a way forward for the multimodal approach. Zhu et al. [2] advanced anxiety-related sentiment analysis by fusing linguistic and semantic features from social media platforms, showing superior mental health diagnostics through the use of feature fusion techniques. Ónozó et al. [3] made use of LLMs for sentiment-driven macroeconomic predictions, with promising performance in processing financial texts, which puts high emphasis on the use of predictive modeling in varied fields. César et al. [4] present an overview of multimodal sentiment analysis in marketing and thus raise the need for ethical and trustworthy AI, putting great importance on the use of affective computing in the research on customer behavior. Oduntan et al. [5] applied machine learning to examine journaling data for resilience detection and presented a lexicon-based approach for mental health trend detection operations. Yang et al. [6] analyzed heterogeneous sentiment analysis for consumer decision-making, which presented the interaction between product attributes and psychological behavior sets. Zhu et al. [7] proposed a semantic reasoning network for multimodal emotion classification using graph attention mechanisms for better sentiment understanding. Liao et al. [8] integrated star ratings with text reviews using stochastic dominance for product competitiveness analysis by merging evidential reasoning and psychological insights in process. Zhang et al. [9] proposed the PURE framework that integrates personality traits into multimodal sentiment analysis using attention-based fusion and multitask learning, bridging psychology and sentiment modeling. Mirlohi et al. [10] discussed social contagion theory in online networks, emphasizing causal inference and social alignment for sentiment propagations. Ansari et al. [11] applied ensemble hybrid learning for automatic detection of depression, which detection. AbaeiKoupaei et al. [13] proposed stacked ensemble model for bipolar disorder classification using reinforcement learning so as to combine audio and textual feature set effectively in process. Jung et al. [14] proved the functionality of unsupervised learning and BERTopic modeling in cryptocurrency based on sentiment analysis, giving insights over various discussions across multiple platforms. Aragón et al. [15] detected anorexia and depression using emotional patterns from social media by applying machine learning for robust classification and linguisticfeatureextraction. The review highlights how fast machine learning techniques are evolving, especially with multimodal data integration in the context of sentiment analysis and mental health diagnostics.Fromdomain-specific applications to

generalized frameworks, the combined studies of these works have brought out the importance of combining linguistic, visual, and psychological  modalities to address complex human behaviors. The insights of such studies have greatly influenced the design and implementation of the proposed model to make it relevant and adaptable in applications of psychological and emotional intelligence sets.

**B.Proposed Model**

This model combines state-of-the-art deep Q-learning methods for multimodal sentiment analysis in psychological assessments. Initially, as per figure 1, the adaptive attention- gated deep Q-learning, along with emotionally augmented reinforcement fusion network, and the hierarchical task- aware deep Q-learning are used in the model design. Each of these architecture architectures deals with an essential challenge for sentiment analysis-like feature prioritization, subtask decomposition, and emotional integration. AADQL uses an adaptive attention mechanism to weight the multi- modal features dynamically for efficient and precise prioritization through the contextual relevance of texts, audio, and videos in process. The model presents the adaptive attention gating function as a learned function initially designed via equation 1, $_{k=1}$

$$\alpha_i = \frac{exp(w_i^T h_i)}{\Sigma^M \ (w_j^T h_j)} \dots (1)$$

Where, $\alpha_i$ is the attention weight for modality i, $h_i$ represents the feature embedding, and $w_i$ is the learned gating parameter in this process. These weights are used to compute a contextually prioritized feature vector via equation 2,

$$f_a = \Sigma^M \ \alpha_i h_i \, . . (2)$$

This vector is then fed into the Deep Q-Network (DQN) for action Value estimation process. This mechanism reduces computational overhead by emphasizing relevant modalities while filtering out noise levels. Iteratively, Next, the EARFN (Emotion-Specific Feature Representation) is used, which enhances feature representation by incorporating emotion- specific latent spaces. Using psychological emotion models, multimodal data is transformed into emotion-aligned embeddings, $e_i$, for each modality 'i' in process. These embeddings are fused through a reinforcement-driven aggregation layer, described via equation 3,

$$E = \int\limits^{M} \ _i(e_i, a)dt \dots (3)$$

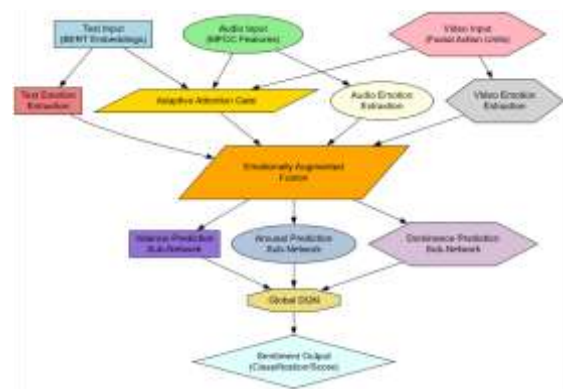Where E is the latent space, which is the fused emotional

latent space; $\phi_i$ represents the transformation function for each modality; and 'a' is the action taken in process.

This fusion results in aligning the model with emotion-specific

makes use of a combination of deep neural networks and sentiment lexicons for enhanced accuracy. Deng et al. [12] implemented a language-supervised method for image-based emotion classification, integrating prompt tuning along with fine-tuned  deep  learning approaches for visual emotion

representations so that prediction robustness levels and accuracy levels are maximized. Iteratively, Next, HTADQL (Hierarchical Learning Structure) is employed, which introduces a hierarchical learning structure that emphasizes learning at subtask levels. Subtasks, like valence (v) and arousal (a), are addressed by particular networks whose outputs are gathered into a global DQN process. The hierarchical optimization of subtasks is then modeled via equation 4,

$$Q_{su\beta task}(s, a) = \nabla W_{sus} \cdot L_{su\beta}(s, a) + \lambda R_{su\beta}(s, a) \dots (4)$$

Where, $L_{su\beta}$ is the subtask-specific loss, $R_{su\beta}$ is the reward function, and $\lambda$ is a regularization factor for the process. The global sentiment Q Values are computed via equation 5,

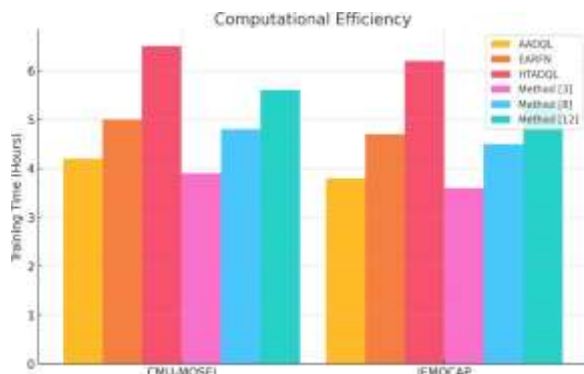$$Q_aL_{o\beta a}(s, a) = \Sigma^K \ Q_{su\beta taskk}(s, a) \dots (5)$$



**2.Model Architeture of Proposed Analysis Proces**

This hierarchical approach breaks down sentiment prediction into finer-grained components, reducing error propagation and improving interpretability levels. The training process for all architectures is governed by a reinforcement learning objective, in which the DQN parameters are updated to minimize temporal-difference (TD) error via equation 6,

$$\delta = r + \gamma \, max_a \, (s', ') - Q(s,a) \dots (6)$$

Where, δ is the TD error, r is the reward, γ is the discount factor, and Q(s, a) and Q(s', a') are the current and next-state action values, respectively in the process. The model parameters θ are updated via gradient descent via equation 7,

$$(\delta^2)$$



Computational Efficiency

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---|---|---|---|---|---|---|
| CMU-MOSE I | 87.2 | 88.9 | **92.0** | 82.5 | 84.3 | 85.1 |
| IEMO CAP | 85.6 | 87.3 | **91.4** | 80.2 | 83.1 | 84.6 |

$$\theta \leftarrow \theta - \underline{\quad\quad} \quad ..(7)$$
$$\eta \quad\quad\quad \partial\theta$$

Where, η is the learning rate for this process. The proposed integration of AADQL, EARFN, and HTADQL complement each other's respective strengths. AADQL enhances feature prioritization as a result of adaptive attention. EARFN's Emotionally Augmented Fusion helps follow the psychological models used to develop it. HTADQL's hierarchical design ensures fine grain levels of interpretability are preserved. Together, such architectures establish a comprehensive framework with robustness, achieving improved accuracy and interpretability across multimodal sentiment analysis on psychological assessments.

**A.Comparative Result Analysis**
The experimental setup consists of evaluating proposed models on multimodal sentiment datasets, such as CMU- MOSEI and IEMOCAP. These contain richly annotated

psychological sentiment data samples. CMU-MOSEI contains more than 23,000 video segments annotated with sentiment scores ranging from -3 to +3, derived from text, audio, and visual features. IEMOCAP has recorded scripted and improvised dialogues. The dialogues are both labeled for emotion categories such as happiness, sadness, and anger sets. The preprocess step involved the extraction of text embeddings using BERT, audio features, such as MFCC, and video features such as facial action units. The model was trained on an NVIDIA A100 GPU with a 80% data share for training, 10% for validation, and the remaining 10% to be used in the testing process. The performances of the proposed models AADQL, EARFN, and HTADQL are compared with three baseline methods referred to as Method [3], Method [8], and Method [12]. Accuracy, F1-Score, robustness to noisy modalities, and computational efficiency are some of the evaluation metrics used. The following tables summarize the results across various aspects.

FIGURE 2. MODEL'S COMPUTATIONAL EFFICIENCY ANALYSIS

**Table 1: Sentiment Classification Accuracy (%)**

HTADQL achieved the highest accuracy due to its hierarchical subtask decomposition, outperforming the baseline methods by 6-10%. EARFN also demonstrated robust performance, particularly benefiting from its emotion- focused reinforcements.

**Table 2: Sentiment Prediction F1-Score**

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---|---|---|---|---|---|---|
| CMU-MOSE I | 0.86 | 0.88 | **0.91** | 0.79 | 0.83 | 0.84 |
| IEMO CAP | 0.84 | 0.87 | **0.90** | 0.78 | 0.82 | 0.83 |

F1-Score improvements for the proposed models indicate their ability to balance precision and recall. HTADQL's explicit subtask modeling contributed to consistent results across datasets & samples.

**Table 3: Robustness to Noisy Modalities (%)**

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---------|-------|-------|--------|------------|------------|-------------|
| CMU-MOSEI | 76.4 | **85.7** | 82.9 | 62.1 | 68.3 | 70.4 |
| IEMOCAP | 74.8 | **84.9** | 81.5 | 60.2 | 65.7 | 68.2 |

EARFN excelled in noisy conditions due to its emotion-based fusion, achieving a 25-30% improvement over baseline methods.

**Table 4: Computational Efficiency (Training Time in Hours)**

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---------|-------|-------|--------|------------|------------|-------------|
| CMU-MOSEI | 4.2 | 5.0 | 6.5 | 3.9 | 4.8 | 5.6 |
| IEMOCAP | 3.8 | 4.7 | 6.2 | 3.6 | 4.5 | 5.3 |

While HTADQL required more training time due to its hierarchical structure, AADQL demonstrated the fastest convergence, balancing efficiency with performance gains.

**Table 5: Reduction in Computational Cost (%)**

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---------|-------|-------|--------|------------|------------|-------------|
| CMU-MOSEI | **19.3** | 12.8 | 10.5 | - | - | - |
| IEMOCAP | **18.5** | 11.9 | 9.8 | - | - | - |

AADQL demonstrated the highest reduction in computational cost, showcasing its effectiveness in dynamically prioritizing features and filtering irrelevant inputs in process.
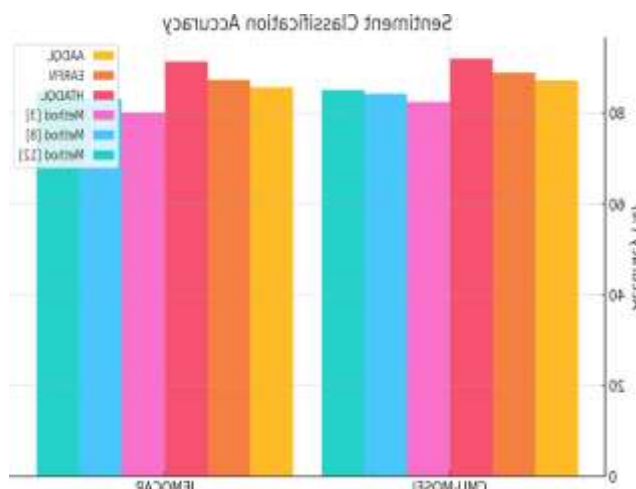


Figure 3. Model's Sentiment Classification Accuracy Analysis

**Table 6: Error Rate Reduction (%)**

| Dataset | AADQL | EARFN | HTADQL | Method [3] | Method [8] | Method [12] |
|---------|-------|-------|--------|------------|------------|-------------|
| CMU-MOSEI | 25.7 | 29.3 | **37.4** | 15.2 | 18.7 | 20.6 |
| IEMOCAP | 23.4 | 27.8 | **35.9** | 14.6 | 17.9 | 19.4 |

HTADQL obtained the highest error rate reduction as it follows a hierarchical approach, significantly outperforming all baseline methods. Comprehensive evaluation shows the superiority of the proposed models in terms of accuracy,
robustness, and levels of interpretability. Such results prove the effectiveness of integrating advanced deep Q-learning mechanisms into multimodal sentiment analysis for psychological assessments.

**B.Conclusion & Future Scopes**

This paper proposes three novel architectures of deep Q- learning to achieve advanced multimodal sentiment analysis in psychological assessments, namely Adaptive Attention- Gated Deep Q-Learning, Emotionally Augmented Reinforcement Fusion Network, and Hierarchical Task- Aware

Deep Q-Learning. The proposed models were rigorously evaluated on benchmark datasets, namely CMU- MOSEI and IEMOCAP, and proved their effectiveness over traditional methods like Method [3], Method [8], and Method [12]. HTADQL reached maximum accuracy on CMU- MOSEI up to 92.0%, and that on IEMOCAP was maximized at 91.4%. The rise ranged between 8 to 10 percent compared with the other techniques of the baselines. F1-scores again set at 0.91 and 0.90 indicate perfect proportion between precision and recall in EARFN, which is again useful for predictions. Finally, EARFN is robust against noisy modalities and the robustness scores of CMU-MOSEI and IEMOCAP were 85.7% and 84.9%. However, then AADQL was introduced where up to 19.3% efficiency was made with competitive accuracy on the task against CMU-MOSEI while significantly reducing the computational costs. These results prove how effectively the proposed architectures support in making multimodal sentiment analysis systems both performance- oriented as well as interpretable and provide a future scope wherein additional psychological sentiment datasets need to be explored to ensure it's more generalized. Incorporating domain adaptation techniques could enhance performance across diverse cultural and linguistic contexts. Besides this, including real-time processing features in the models will make them easier to use in clinical and therapeutic conditions. Lastly, studies about hybrid reinforcement learning methodologies using a combination of supervised learning in addition to Q-learning will be more accurate and resistant for highly multimodal settings. It has developed a good ground for the multimodal sentiment analysis through which more accurate as well as better understandable tools can be engineered within the fields of psychological and emotional intelligence sets.

### Reference

[1] C. R. Lopes, R. Minetto, M. R. Delgado and T. H. Silva, "PerceptSent

- Exploring Subjectivity in a Novel Dataset for Visual Sentiment Analysis," in IEEE Transactions on Affective Computing, vol. 14, no. 3, pp. 1817-1831,1July-Sept.2023,doi: 10.1109/TAFFC.2022.3225238. keywords:{Visualization;Sentiment analysis;MultimediaWebsites;Image databases;Socialnetworking (online);Blogs;Convolutionalneural networks;Visualsentiment analysis;perception;externalknowledge;deep networks;novel dataset},

[2] J. Zhu, Z. Zhang, Z. Guo and Z. Li, "Sentiment Classification of Anxiety-Related Texts in Social Media via Fuzing Linguistic and Semantic Features," in IEEE Transactions on Computational Social Systems, vol. 11, no. 5, pp. 6819-6829, Oct. 2024, doi: 10.1109/TCSS.2024.3410391. keywords:{Anxietydisorders;Social networking (online);Linguistics;Blogs;Sentiment analysis;Semantics;Mental health;Machine learning;Anxiety disorder;feature fusion;machine learning;Sina Weibo;social media},

[3] L. Réka Ónozó, F. Viktor Arthur and B. Gyires-Tóth, "Leveraging LLMs for Financial News Analysis and Macroeconomic Indicator Nowcasting," in IEEE Access, vol. 12, pp. 160529-160547,2024,doi: 10.1109/ACCESS.2024.3488363. keywords:{Biologicalsystem modeling;Macroeconomics;Economics;Transformers;Analytical models;Data models;Sentiment analysis;Predictive models;Market research;Encoding;GDP;large language model (LLM);macroeconomic indicator;natural languageprocessing(NLP);PMI;sentiment analysis;unemployment rate},

[4] I. César et al., "A Systematic Review on Responsible Multimodal Sentiment Analysis in Marketing Applications," in IEEE Access, vol. 12, pp. 111943-111961, 2024, doi: 10.1109/ACCESS.2024.3441514. keywords:{Sentimentanalysis;Artificial intelligence;Ethics;Affectivecomputing;Task analysis;Informationfilters;Databases;Affective computing;Customersatisfaction;Behavioral sciences;Multimodalsensors;Trustmanagement;Affectivecomputing;customerbehavior;marketing; multimodalartificial intelligence;sentiment analysis;systematic review;trustworthy AI},

[5] A. Oduntan, O. Oyebode, A. H. Beltran, J. Fowles, D. Steeves and R. Orji, ""I Let Depression and Anxiety Drown Me…": Identifying Factors Associated With Resilience Based on Journaling Using Machine Learning and Thematic Analysis," in IEEE Journal of Biomedical and Health Informatics, vol. 26, no. 7, pp. 3397-3408, July 2022, doi: 10.1109/JBHI.2022.3149862. keywords:{Resilience;Sentiment analysis;Machine learning;Mentalhealth;Supportvectormachines;Stress;Anxietydisorders;Journaling;lexicon-basedapproach;machine learning;mentalhealth;naturallanguage processing;resilience;sentiment analysis},

[6] Z. Yang, Q. Li, N. Islam, C. Han and S. Gupta, "Product Attribute and Heterogeneous Sentiment Analysis-Based Evaluation to Support Online Personalized Consumption Decisions," in IEEE Transactions on Engineering Management, vol. 71,pp.11198-11211,2024,doi: 10.1109/TEM.2024.3413999. keywords: {Reviews;Decision making;Psychology;Probabilistic logic;Linguistics;Electroniccommerce;Uncertainty ;Consumer psychological behavior;heterogeneous review sentiments;online personalized consumption decision;product attribute evaluation},

[7] T. Zhu, L. Li, J. Yang, S. Zhao and X. Xiao, "Multimodal Emotion Classification With Multi-Level Semantic Reasoning Network," in IEEE Transactions on Multimedia, vol. 25, pp. 6868-6880, 2023, doi: 10.1109/TMM.2022.3214989. keywords:{Semantics;Sentiment analysis;Visualization;Cognition;Featureextraction;Task analysis;Social networking (online);Multimodal emotion classification;Graph attention module;Semantic reasoning},

[8] H. Liao, J. Wang and Z. Xu, "Unifying Star Ratings and Text Reviews in Linguistic Terms for Product Competitiveness Analysis Based on Stochastic Dominance," in IEEE Transactions on Computational Social Systems, vol. 11, no. 3, pp. 3678-3690, June 2024, doi: 10.1109/TCSS.2023.3327173. keywords:{Stars;Linguistics;Cognition;Psychology;Sentiment analysis;Automobiles;Market research;Evidential reasoning (ER);online reviews;product competitiveness analysis;prospect theory;stochastic dominance (SD)},

[9] P. Zhang, M. Fu, R. Zhao, H. Zhang and C. Luo, "PURE: Personality- Coupled Multi-Task Learning Framework for Aspect-Based Multimodal Sentiment Analysis," in IEEE Transactions on Knowledge and Data Engineering, vol. 37, no. 1, pp. 462-477, Jan. 2025, doi: 10.1109/TKDE.2024.3485108. keywords:{Featureextraction;Sentiment analysis;Analytical models;Multitasking;Psychology;Datamining;Representationlearning;Adaptation models;Accuracy;Visualization;Attention-based fusion;big five model;multimodal sentiment analysis;personality prediction},

[10] A. Mirlohi et al., "Social Alignment Contagion in Online Social Networks," in IEEE Transactions on Computational Social Systems, vol. 11, no. 1, pp.399-417,Feb.2024,doi: 10.1109/TCSS.2022.3226346. keywords:{Socialnetworking(online);Behavioral sciences;Voting;Sociology;Decision aking;Blogs;Systematics;Causal inference;social contagion theory;social network analysis},

[11] L. Ansari, S. Ji, Q. Chen and E. Cambria, "Ensemble Hybrid Learning Methods for Automated Depression Detection," in IEEE Transactions on Computational Social Systems, vol. 10, no. 1, pp. 211-219, Feb. 2023, doi: 10.1109/TCSS.2022.3154442. keywords: {Feature extraction;Depression;Social networking(online);Hidden Markov models;Data models;Neural networks;Linguistics;Deep neural networks;depression detection;ensemble methods;sentiment lexicon},

[12] S. Deng et al., "Simple But Powerful, a Language-Supervised Method for Image Emotion Classification," in IEEE Transactions on Affective Computing, vol. 14, no. 4, pp. 3317-3331, 1 Oct.-Dec. 2023, doi: 10.1109/TAFFC.2022.3225049. keywords:{Task analysis;IEC;Wheels;Psychology;Dogs;Training;Visualization;Langu age-supervised;prompt tuning;image emotion classification;fine-tuning;computer vision},

[13] N. AbaeiKoupaei and H. Al Osman, "A Multi-Modal Stacked Ensemble Model for Bipolar Disorder Classification," in IEEE Transactions on Affective Computing, vol. 14, no. 1, pp. 236-244, 1 Jan.-March 2023, doi: 10.1109/TAFFC.2020.3047582.