# M²Lens: Visualizing and Explaining Multimodal Models for Sentiment Analysis

Xingbo Wang, Jianben He, Zhihua Jin, Muqiao Yang, Yong Wang, and Huamin Qu
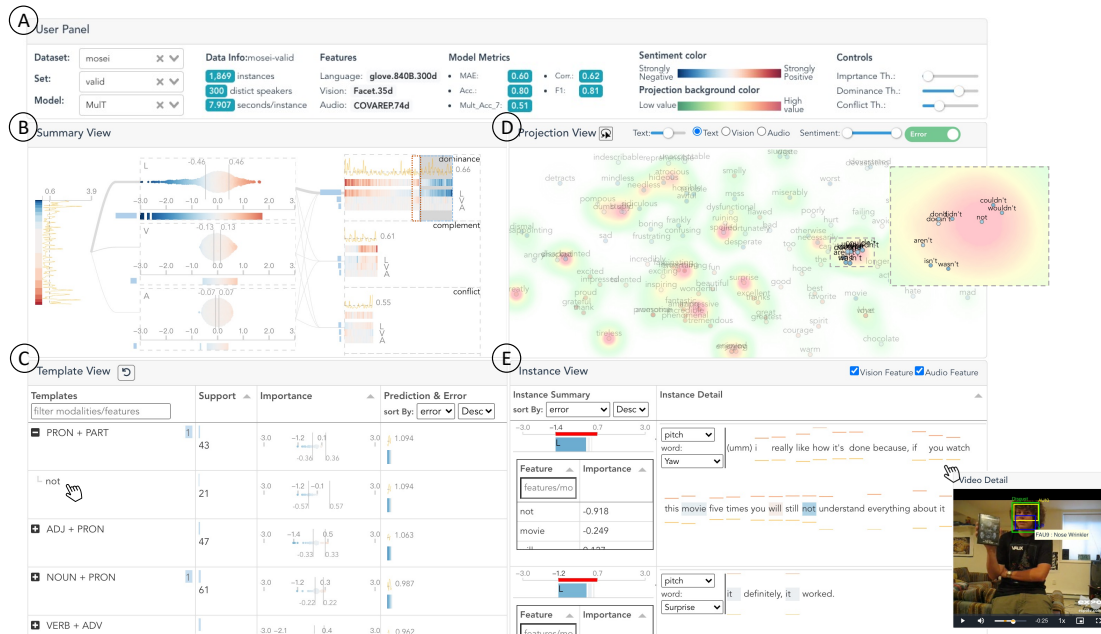
Fig. 1. The explanatory interface of *M²Lens* consists of five views. The *User Panel* (A) displays the descriptive statistics about the model and dataset. The *Summary View* (B) presents a global summary of the importance of individual modalities, as well as their interactions using a three-layer augmented tree-like layout. The *Template View* (C) and *Projection View* (D) complement each other for subset-level explanations. Specifically, *Template View* (C) summarizes frequent and influential templates of feature sets in a table. The *Projection View* (D) supports multi-faceted explorations of instances that have features of interest. The *Instance View* (E) provides local explanations by visualizing the important features and the context of individual instances.

**Abstract**—Multimodal sentiment analysis aims to recognize people's attitudes from multiple communication channels such as verbal content (i.e., text), voice, and facial expressions. It has become a vibrant and important research topic in natural language processing. Much research focuses on modeling the complex intra- and inter-modal interactions between different communication channels. However, current multimodal models with strong performance are often deep-learning-based techniques and work like black boxes. It is not clear how models utilize multimodal information for sentiment predictions. Despite recent advances in techniques for enhancing the explainability of machine learning models, they often target unimodal scenarios (e.g., images, sentences), and little research has been done on explaining multimodal models. In this paper, we present an interactive visual analytics system, *M²Lens*, to visualize and explain multimodal models for sentiment analysis. *M²Lens* provides explanations on intra- and inter-modal interactions at the global, subset, and local levels. Specifically, it summarizes the influence of three typical interaction types (i.e., dominance, complement, and conflict) on the model predictions. Moreover, *M²Lens* identifies frequent and influential multimodal features and supports the multi-faceted exploration of model behaviors from language, acoustic, and visual modalities. Through two case studies and expert interviews, we demonstrate our system can help users gain deep insights into the multimodal models for sentiment analysis.

**Index Terms**—Multimodal models, sentiment analysis, explainable machine learning

◆

## 1 INTRODUCTION

- *Xingbo Wang, Jianben He, Zhihua Jin, and Huamin Qu are with the Hong Kong University of Science and Technology. E-mail: {xingbo.wang, jhebt, zjinak, huamin}@ust.hk.*
- *Muqiao Yang is with Carnegie Mellon University. E-mail: muqiaoy@andrew.cmu.edu.*
- *Yong Wang is with Singapore Management University. E-mail: yongwang@smu.edu.sg. He is the corresponding author.*

Sentiment analysis aims to use computational approaches to identify people's attitudes, opinions, and other subjective information in human communication. It can benefit various applications, such as customer analysis, social robots, and political campaigns. Prior research on sentiment analysis is mainly based on a single communication channel (i.e., text or facial expression) [60, 64, 82], which is often referred to as *unimodal sentiment analysis*. However, human communication is often multimodal. For example, people can show their happiness through positive words and tones, along with a wild smile. With the thriving of social media, a large number of multimodal communication datasets

can be collected and studied, e.g., TV series and vlogs showing people's sentiment towards different topics and objects. This has greatly boosted the development of *multimodal sentiment analysis* techniques, and it has already become a vibrant and important research topic.

Unlike the long-established unimodal sentiment analysis, multimodal sentiment analysis combines the heterogeneous data and captures two primary forms of interactions in different modalities: ***intra-modal*** and ***inter-modal*** interactions. Intra-modal interactions refer to the dynamics of one modality, which is the same as the unimodal analysis based on the single communication channel. Inter-modal interactions consider the correspondence between different modalities across time, e.g., the co-occurrences of a happy tone and a smile or a sudden pause after a humorous punchline. In practice, people's communication styles are highly complex and idiosyncratic. For example, a sentence may seem semantically positive, but people can express it with a sarcastic tone to reveal their dissatisfaction. In such cases, unimodal sentiment analysis is not reliable, while multimodal models can offer the opportunities to explore vocal and visual expressions besides texts. In addition, previous research [4, 47, 64] has confirmed that multimodal models are more accurate and robust in various downstream tasks.

Currently, deep-learning-based models achieve superior performance over the traditional methods [40, 75] in multimodal sentiment analysis. Representative examples include transformers [50, 69], Convolutional Neural Networks (CNNs) [72], and Recurrent Neural Networks (RNNs) [48, 51, 78]. However, these models often work like black-boxes, hindering users from understanding the underlying model mechanism and fully trusting them in decision-making. Enhancing the explainability of deep learning models has become critical for both model developers and users, and received increasing attention in the past few years [3, 23]. For example, post-hoc explainability techniques, such as LIME [54], SHAP [38], and IG [67], help identify important features (e.g., words or image patches) that influence model predictions. However, these methods often target providing local explanations on instances (e.g., sentences) in unimodal scenarios. They do not scale well to produce global explanations on how intra- and inter-modal interactions influence the model decisions, for example, how the models will behave when positive words and sad voices are presented.

It is challenging to explain multimodal models for sentiment analysis. First, it is necessary to relate the model performance back to the multimodal input data [44, 53]. The heterogeneity and high dimensionality of multimodal human behaviors make it difficult for users to easily interpret the input features or data, as well as how they affect model decisions. Compact and human-friendly summaries of multimodal data are highly desired, but little research (if not no) has been done on it. Second, it is non-trivial to explain inter-modal interactions between different modalities explicitly, which, however, are the unique characteristics of multimodal sentiment analysis models. For example, when a person says something positive with a neutral voice and facial emotion, users may feel interested in whether the models can discern positive sentiment in the language modality (i.e., the text).

In this paper, we propose $M^2Lens$, a novel explanatory visual analytics tool to help both developers and users of multimodal machine learning models better understand and diagnose Multimodal Models for sentiment analysis. By considering the feature importance measured by post-hoc explainability techniques, $M^2Lens$ interprets intra- and inter-modal interactions learned by a multimodal language model from the global, subsets, and local levels. Particularly, we focus on interpreting three typical types of interactions, i.e., ***dominance***, ***complement***, and ***conflict***. Moreover, it facilitates a multi-faceted exploration of the multimodal features and their influences on the model decisions for sentiment analysis. $M^2Lens$ consists of four major views. Specifically, the *Summary View* features an augmented tree-like layout for global explanations of the impacts of individual modalities and their interplay. The *Template View* summarizes influential and frequent multimodal features with compact templates. The *Projection View* enables multi-faceted exploration of the features of user interest using different glyph designs. The *Instance View* visualizes individual multimodal instances and their explanations with details.

In summary, our major contributions are:

- $M^2Lens$, a visual analytics system to produce multi-level and multi-faceted explanations on intra- and inter-modal interactions that are learned by the multimodal models.
- Case studies and expert interviews that demonstrate the effectiveness of our approach in helping users gain deep insights into two state-of-the-art multimodal models for sentiment analysis.

## 2 RELATED WORK

This section discusses the relevant research of our approach, including multimodal language analysis, post-hoc explainability techniques, and machine learning interpretation with visualization.

### 2.1 Multimodal Sentiment Analysis

Multimodal sentiment analysis is a vibrant topic in natural language processing (NLP). It automatically extract people's attitudes or affective states from multiple communication channels (e.g., text, voice, and facial expressions). Moreover, it has various applications [24, 80, 81]. The core challenge is modeling the complex *intra-modal* and *inter-modal* interactions, where multimodal features are being fused.

Early work [35, 41] concatenated features from different modalities before being input to a learning model. Conversely, some work adopted *late-fusion* approaches that combine the decision values from individual unimodal models using a voting scheme [42, 49] or a learning model [17, 52]. However, these methods ignore the cross-modal interactions. To address such issues, some work explicitly computed the unimodal, bimodal, and trimodal features and fused them with tensor product [36, 76] and dynamic routing [70]. Recently, neural network methods [11, 46, 51, 69, 77, 78] are popular to model the complex interplay between modalities. For example, researchers [11, 51] have extended LSTM cells and gates to learn temporal interaction patterns among multimodal sequences. Pham et al. [46] proposed attention-based RNNs to learn multimodal representations with a cyclic translation loss among modalities. Zadeh et al. [77] designed a multi-view gated memory unit that is controlled by neural networks. It stores and predicts temporal cross-modal interactions. Tsai et al. [69] utilized transformer attention mechanisms to learn both cross-modal alignment and interactions. Although neural networks greatly improve the performance over traditional methods, their complex architecture seriously affects the model interpretability. This paper presents an explanatory interface to diagnose black-box models for sentiment analysis tasks.

### 2.2 Post-hoc Explainability Techniques

Post-hoc explainability techniques interpret models after the training process [3, 37]. They generally include *model-specific* and *model-agnostic* approaches [37]. Model-specific methods explain particular models ranging from shallow models [18, 68] to sophisticated neural networks [30, 61]. In contrast, model-agnostic methods are flexible enough to be applied to any machine learning model. Here, we discuss two main types of model-agnostic approaches: explanation by simplification and feature relevance explanation [3].

For simplification techniques, researchers often built surrogate models (e.g., rule-based learners [25, 29, 55], decision trees [5], and linear models [54]) to imitate the original model behaviors with reduced complexity. One of the most representative methods is LIME [54], which builds locally linear models to approximate individual predictions based on neighbors of instances of interest. Feature relevance explanation quantifies the feature contributions to model predictions. One popular example is SHAP [38], whose mathematical root is Shapley Value [62]—a method from cooperative game theory. SHAP computes an additive importance score for each feature to describe its influence, given a prediction result. It has desirable properties (local accuracy, missingness, and consistency) and is proved to be aligned with human intuitions. Other work used local gradients [57], randomized feature permutations [21], or influence functions [28] to disclose feature relevance.

However, the methods above are often used to interpret specific instances of one modality (e.g., sentences, images), which cannot be directly applied to multimodal sentiment analysis. This paper aims to fill the gap by enabling multi-level explanations on the learned intra- and inter-modal interactions from global, subsets, and local levels.

## 2.3 Machine Learning Interpretation With Visualization

With the increasing complexity of both data and machine learning models, various visual analytics systems have been proposed to assist in understanding the model behaviors. Besides measuring the model performance with computational metrics, users also need to explore when and why a model makes specific decisions [23]. One of the most common and important interpretation strategies in previous work is to reveal the relationship between the *input data* and *model predictions* [3, 23]. They can be categorized into two groups: *instance exploration* and *feature & subset exploration*.

Instance visualization shows model behavior towards individual data samples. Amershi et al. [2] presented ModelTracker to support performance debugging with a visual summary of binary classification instances. Ren et al. [53] extended the performance visualization to multi-class scenarios with aligned vertical axis designs, while Kahng et al. [26] and Alsallakh et al. [6] adopted a matrix-like design for instance summary. Apart from visualizing instance distributions, Kulesza et al. [34] built an exploratory debugging prototype to enable users to explain corrections back to models. In addition, there are tools [20, 63] that allow users to interactively probe models with provided inputs.

Feature and subset visualization investigates how to surface the patterns groups of features [7, 31, 32] and instances [1, 9, 74, 83] that affect model decisions. Brooks et al. [7] developed FeatureInsight, which supports the feature ideation process with a visual summary of set errors. Krause et al. [31] enabled exploration of the predictive power of feature candidates across different feature selection algorithms. For specific applications in CV and NLP, features are often visualized as image patches [43, 61, 65] or text segments [16, 27]. Besides, researchers built interactive tools to facilitate group-level exploration. Zhang et al. [83] conducted feature attribution comparisons to inspect discrepancies across different data subsets. Some work [1, 9, 74] used fairness metrics to partition data into groups for model diagnosis.

However, these methods do not consider exploring multimodal features and determining how much they affect model decisions. Our system facilitates multi-faceted exploration of multimodal features and generates multi-level visual explanations on their influences.

## 3 BACKGROUND

In multimodal sentiment analysis, a machine learning model predicts sentiment based on the visual, acoustic, and language features extracted from the raw video data. This section introduces the related background about multimodal datasets, feature engineering techniques, performance metrics, and intra- and inter-modal interactions.

## 3.1 Dataset

There is a wide range of multimodal datasets in the community. For example, **IEMOCAP** [8] contains 151 videos of dialogues with different emotion labels. **YouTube** [40] consists of videos of product reviews extracted from the social media website, YouTube. Without loss of generality, our work focuses on the largest and widely-used benchmark dataset for multimodal sentiment analysis, i.e., **CMU-MOSEI** [79]. It consists of 23,454 monologue movie review video clips from 1,000 speakers and 250 topics in YouTube. The sentiment of each video clip is labeled by three annotators with a Likert scale of $[-3, 3]$, where 3 indicates strongly positive, $-3$ represents strongly negative, and 0 means neural. Besides the sentiment label, each video is associated with the information from the three communicative channels—transcripts for language resources ($l$), facial expressions for the visual ($v$), and voice of speakers as the acoustic modalities ($a$).

## 3.2 Multimodal Feature Engineering

Prior research on multimodal models mostly uses different feature engineering techniques for all three modalities in sentiment analysis. Here, we follow the common practice of multimodal feature extraction (also provided by CMU-MOSEI). For language features, transcripts are encoded by high-dimensional word vectors. We leverage Glove embeddings [45] to represent each word, where each word is transformed to a 300-dimension vector. For visual modality, most work focuses on facial expressions, which are often encoded by Facial Action Coding System

(FACS) [15]. FACS encodes the facial muscle movement with 35 facial action units. We deploy it to extract frame-level facial features. The acoustic features are engineered through a speech processing framework, COVAREP [13]. The extracted features have 74 dimensions, and all of them are related to speech emotions and tones. To help users gain a quick overview of these fundamental features, we further group them into different classes, which will be introduced in Sect. 5.2.2.

## 3.3 Metrics for Multimodal Sentiment Analysis

Prior work applies several metrics to evaluate the model performance for multimodal sentiment analysis, including mean absolute error ($MAE$), the correlation between the model predictions and human labels ($Corr.$), F1 score ($F1$), 7-class accuracy ($Acc_7$), and 2-class accuracy ($Acc$). Note that $Acc_7$ considers all of the sentiment scores $\mathbb{Z} \in [-3, 3]$, while $Acc$ is a binary classification score that only predicts whether this video clip is positive or negative.

## 3.4 Intra- and Inter-modal Interactions

In practice, sentiment analysis relies on multimodal language signals (e.g., language, facial expressions, and tones). A successful multimodal sentiment analysis requires the understanding of the combinations of these signals, where two primary forms of interactions exist—***intra- and inter-modal*** interactions [4, 64].

When modeling intra- and inter-modal interactions, three typical situations arise [4, 64, 77]:
- One modality is ***dominant*** for sentiment analysis. For example, people may show agreement by nodding their heads, where the vision modality dominantly indicates their positive attitudes.
- More than one modalities ***complement*** each other when people are expressing their sentiment. For example, people's positive attitudes in words can be enhanced by a happy tone.
- More than one modalities ***conflict*** with each other. For example, people may tell sad stories with smiles on their face.

Researchers have tried to build models to analyze the situations above for better sentiment analysis. However, most state-of-the-art models are deep-learning-based techniques with little interpretability. Model developers and users are not aware of how exactly the model utilizes information in multiple modalities in situations of dominance, complement, or conflict. Explaining multimodal model behaviors not only provides insights into the multimodal language characteristics, but also reveals the model errors and inspires new model designs. In our work, we explicitly provide global explanations on intra- and inter-modal interactions with a compact visual summary. Specifically, we categorize instances into dominance, complement, and conflict groups based on the importance of each modality computed by SHAP [38]. Furthermore, we summarize influential feature sets for each group with templates to provide finer-grained explanations on model behavior.

## 4 DESIGN REQUIREMENTS

Our goal is to develop a visual analytics system to help users (e.g., model developers and model users) understand and diagnose the behaviors of multimodal models for sentiment analysis. Similar to the general black-box explanation tools [2, 32, 53, 83], interpreting multimodal models helps target users gain insights into the connection between the model performance (e.g., model errors) and the characteristics of multimodal data. For example, model users can examine whether a model has a bias or poor performance on some types of data and further decide if it is a proper fit for target applications. Furthermore, given the critical aspects of multimodal sentiment analysis (in Sect. 3.4), it is beneficial to explain the intra- and inter-modal relationships learned by the model. For instance, model developers can adjust the fusion weights of different modalities based on their relative importance to achieve better sentiment predictions. However, it is challenging to interpret multimodal models due to the high complexity of multimodal data and inter-modal relationships.

To understand users' general needs and formulate design requirements, we surveyed prior visualization techniques for interpreting machine learning models [2, 3, 7, 10, 26, 31, 32, 39, 53, 83] and multimodal language analysis [4, 41, 69, 70, 76, 79]. Also, we worked closely with a

researcher in NLP and multimodal machine learning (who is also a co-author of this paper) for about five months to collect his feedback and iteratively refine the design requirements. We summarize the design requirements as follows.

**R1: Show the model performance.** Performance metrics are crucial for guiding the model analysis [2,53]. They provide quantitative measures of how accurate the predictions are and can help users pinpoint where the model is likely to fail. The users often want to evaluate models at different levels:

*Q1: What are the overall error distributions for model predictions?*
*Q2: What are the instances that are predicted with large/small errors?*

**R2: Reveal the contributions of modalities to the model predictions.** Besides performance metrics, the system should provide global explanations on how the model generally works, especially when working with huge datasets [3,10,26,39]. In multimodal sentiment analysis, intra- and inter-modal interactions are crucial for understanding the model behaviors [4,41]. Thus, it is essential to summarize the influences of individual modalities and their interplay for predictions. Specifically, the system should help users answer the following questions:

*Q3: How does each modality influence the model predictions?* Displaying the contributions of each modality helps users prioritize their efforts in diagnosing a particular modality for model predictions [76].
*Q4: Which modalities **dominate** the model predictions? Also, which modalities **complement** or **conflict** with each other for model predictions?* To better reveal the characteristics of multimodal interactions captured by the model, the system should further summarize the instances according to the interaction types [69,70,79]. Specifically, dominant, complementary, and conflicting modalities, which depict typical interaction types, are the targets for analysis.
*Q5: How do **dominant/complementary/conflicting** modalities influence the model predictions?* Besides recognizing the learned interaction types, it is also essential to connect them to the model predictions for a comprehensive understanding of model behaviors [2,53,83]. For example, the dominance of language modality can contribute to positive or negative sentiment for different instances.

**R3: Identify the influences of multimodal features for the model predictions.** With a global understanding of how the model work on individual modality (**R2**), users need to drill down to finer-level inspection on model behaviors. Feature-based exploration is a common and effective approach for explaining machine learning models [7, 31, 32]. Accordingly, the system should connect high-level modality interactions with the corresponding multimodal features. For example, users may want to know when the language modality dominates the predictions and what words people use to express their sentiments.

*Q6: What are the feature sets that significantly contribute to positive/negative sentiment predictions?* Exploring all the features of instances individually is tedious given the high volume and dimensionality of multimodal data. Summarizing the set of features with a significant predictive contribution helps reduce the efforts in exploration [7, 31]. In addition, it helps users develop a high-level concept about model predictions. For example, users may want to know what types of words or facial expressions are considered important to models when dealing with positive sentiment cases.
*Q7: What features are considered important by the model? Are they plausible for prediction?* To help users analyze the individual predictions, features with a significant influence on the model performance should be presented to users and allow them to judge whether they align well with the observation of the original data.

**R4: Support multi-level and multi-faceted exploration of the multimodal model behaviors.** Given the multimodal settings of sentiment analysis, the visualization should empower users to explore the relationships between the model and input data from multiple aspects (e.g., language, facial expressions). To facilitate a comprehensive understanding of multimodal models, explanations should be offered on different levels, including the influences of individual modalities and their interplay, and the importance of multimodal features.

## 5 M²LENS

Based on the derived design requirements (Sect. 4), we develop a visual analytics system, $M^2Lens$ (Fig. 1), for understanding and diagnosing how models utilize multimodal information for sentiment prediction. In this section, we first provide an overview of the system architecture. Then, we will illustrate the methods for generating explanations of multimodal model behavior. Next, we describe the visual designs and interactions in detail.

### 5.1 System Overview

Fig. 2 shows the system architecture. First, speakers' opinion videos are transformed into visual, acoustic, and language features. The *storage* module saves users' model and data with processed features. Then, the *explanation engine* inputs the features into the model and generates multi-level explanations of model behaviors based on the feature attribution methods (e.g., SHAP). The *visual analysis* module enables interactive exploration of the explanations through five main views.
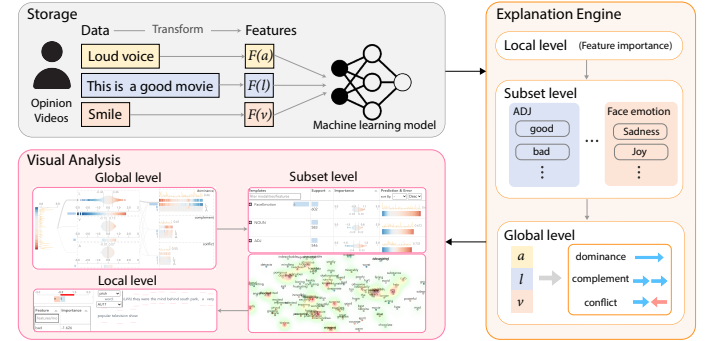


Fig. 2. $M^2Lens$ consists of a *storage* module, an *explanation engine*, and a *visual analysis* interface.

The *User Panel* is the entry point of the whole interface, where the descriptive statistics about the model performance and dataset (***Q1***) are shown. Then, *Summary View*, *Template View*, *Projection View*, and *Instance View* provide multi-level model explanations from language, visual, and acoustic modalities (***R4***). The *Summary View* presents a global summary of the influences of individual modalities and their interplay for the sentiment predictions (***R2***). The *Template View* and *Projection View* complement each other for subset-level explanations (***R3***). Specifically, *Template View* uses templates to summarize feature sets that frequently and significantly contribute to the model predictions. *Projection View* supports the multi-faceted exploration of instances that have features of interest, along with their prediction errors. The *Instance View* summarizes instance-level prediction information (e.g., errors) (***Q2***) and offers local explanations on the importance of each modality and its features (***R4***). In addition, it adds the audio and vision features along the spoken words and provides the corresponding raw video clips with feature annotations for further exploration.

### 5.2 Multi-level Explanations

To facilitate users with a comprehensive understanding of multimodal behavior, we propose methods to generate global and subset-level explanations (***R2, R3***). They supplement the local explanations computed by feature attribution methods.

#### 5.2.1 Global Explanations

Since the intra- and inter-modal interactions lie at the heart of multimodal sentiment analysis, they are essential for users to understand how the multimodal model utilizes the information from different modalities (i.e., language, audio, and vision) (***R2***). In our work, we characterize three typical types of interactions among modalities—dominance, complement, and conflict (details are in Sect. 3.4).

The ***dominance*** suggests that the influence of one modality dominates the polarity (i.e., positive or negative) of a sentiment prediction. The ***complement*** indicates that two or all three modalities affect a

model prediction in the same direction (i.e., positively or negatively). Conversely, the **conflict** reveals that the influences of modalities differ from each other. According to the definitions above, we formulate a set of rules to identify them (Algorithm 1). Specifically, The influence of the interactions on the model output is based on the importance of each modality $(I_l, I_a, I_v)$, which is the summation of the importance of all its features. Then, we extract and summarize the interactions ($L$) with strong influences for all the predictions. The thresholds for our rules are determined by maximizing the distances between the interaction types while minimizing the average influences of interactions that do not belong to dominance, complement, or conflict (i.e., others):

$$\underset{\{Th_{sig}, Th_{dom}, Th_{confl}\}}{\arg\max} \frac{1}{|L|^2} \sum_i^L \sum_j^L dist(L_i, L_j) - \bar{L}_{others} \quad (1)$$

where $L_i$ ($i \in \{\textbf{dominance, conflict, complement, others}\}$) is the interaction types output by Algorithm 1 for all the instances, *dist* is the Euclidean distance between the average influences of $L_i$ and $L_j$.

---

**Algorithm 1** Rules for extracting important relationships of modalities.

**Input:** $\{I_l, I_a, I_v\}$; $Th_{sig}, Th_{dom}, Th_{confl} (\in (0,1))$;
**Output:** Label for the interaction types, $l$;
1: **if** $\forall i \in \{l, a, v\}, |I_i| > Th_{sig}$ **then**
2:     `/* important interactions */`
3:     **if** $\exists i, j \in \{l, a, v\}, I_i \cdot \sum I_j > 0, \frac{|I_i|}{\|I\|} \geq Th_{dom}$ **then**
4:         $l = \textbf{dominance};$
5:     **else if** $\exists i, j \in \{l, a, v\}, I_i \cdot I_j < 0, \sum \frac{I_i}{\|I\|} \leq Th_{confl}$ **then**
6:         $l = \textbf{conflict};$
7:     **else if** $\exists i, j \in \{l, a, v\}, I_i \cdot I_j > 0$ **then**
8:         $l = \textbf{complement};$
9:     **else**
10:         $l = \textbf{others};$
11: **else**
12:     $l = \textbf{others};$

---

### 5.2.2 Feature Templates

Compared with inspecting the impacts of individual features, exploring feature groups is more effective for analyzing complex model behaviors and data characteristics [26, 33]. It helps users develop a mental model about the model decisions (**Q6**). For example, what types of words (e.g., adjectives) are considered important indicators for positive sentiment. To ease the exploration of influences of high-dimensional features, we organize the model's input features introduced in Sect. 3.2 into several meaningful groups. Then, we summarize frequent and influential groups with compact templates (Fig. 1C).

To promote the understanding of model behaviors, we first identify several *feature sets* based on the sentence structures for the language modality, emotion-related features for the acoustic modality, and facial expressions for the visual modality:

- *Language*: part of speech (POS)[1] (e.g., noun, adjective, verb);
- *Audio*: pitch, amplitude, glottal/voice quality, and phase;
- *Vision* (i.e., Face): face parts (i.e., brow, eye, nose, lip, and chin), head movement, and face emotions.

For language modality, POS features provide a compact summary of the structure of language use. They have been widely used as a probe for natural language [56, 58, 66]. The audio features are grouped according to a state-of-the-art speech processing framework, COVAREP [13], and speech applications [59, 73]. These sets generally relate to the emotions and tones of speech. For face-related features, we divide them into the face parts, head movement, and face emotions. They are the representative components in the facial action coding system (FACS) [14] for describing facial expressions. For the mapping between low-level multimodal features and the feature sets, please refer to the supplementary material.

---

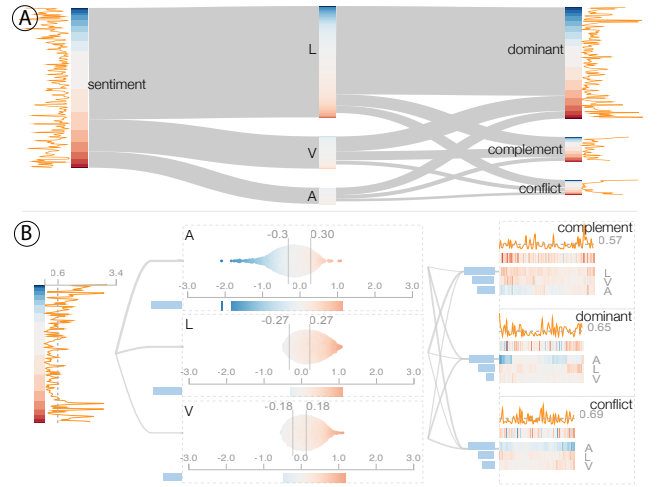[1] https://universaldependencies.org/docs/u/pos/

---

Fig. 3. Design choices for the *Summary View*. A: An augmented Sankey diagram. B: Our current design of augmented tree-like layout.

After grouping the low-level features for each modality, we construct templates for both the frequent feature sets (e.g., "ADJ") and features (e.g., word "good") that have a strong influence on predictions (**Q6**). Specifically, we create itemsets of important features and feature sets for all predictions. Then, we build FP trees [19] to find frequent patterns within the itemsets. For example, if "PRON" and "PART" or the word "not" constantly appear, they will be recorded in the templates (Fig. 1C).

### 5.3 User Interface

Based on the generated explanations, the user interface of $M^2 Lens$ facilitates multi-level exploration of model behavior from the perspective of language, acoustic, and visual modalities (**R4**). All the views are tightly integrated with interactions to ensure a smooth transition between different levels of explanations. They share the same color encoding scheme where dark red means strong positive sentiment and dark blue represents strong negative sentiment.

#### 5.3.1 Summary View

The *Summary View* presents an overview of the intra- and inter-modal interactions that are learned by the selected model in the *User Panel* (**R2**). The influences of individual modalities and their interplay are visualized in a three-layer augmented tree-like layout (Fig. 3B).

**Visual designs.** In the parent node, a barcode chart and a line chart show the distributions of the ground truths and model prediction errors, respectively (**Q1**). The vertical height of the barcode represents the total number of instances, and the color displays the sentiment. Meanwhile, the horizontal position of the line chart suggests the absolute error, and the mean error is represented as a dashed line.

The second layer presents the importance of individual modalities in bee swarm plots (**Q3**). They are arranged according to the influences of modalities in descending order. For each node in the layer, a blue bar is put to the left, whose horizontal length summarizes the total influences of the modality. Besides, the dots in the bee swarm plot and their projections (i.e., the barcode below) demonstrate the distribution of the influences of that modality for all the instances. The color and horizontal position of the dots encode the importance values, while the two gray lines indicate the magnitude of mean absolute importance.

The last layer summarizes the information about the four types of interactions (Sect. 5.2), where the most influential one is shown at the top (**Q4, Q5**). For each interaction, the horizontal range of all its charts marks the number of instances in that group. To better surface the patterns of how the combinations of modalities affect the model predictions, we put the data instances close to each other if all their three modalities share similar influence patterns. Specifically, the similarity is measured by the farthest distances among three modalities between the instances. Then, a line chart and a barcode chart at the top summarize the error and prediction patterns, which are similar to the parent node. In addition, three barcode charts are attached below to present the
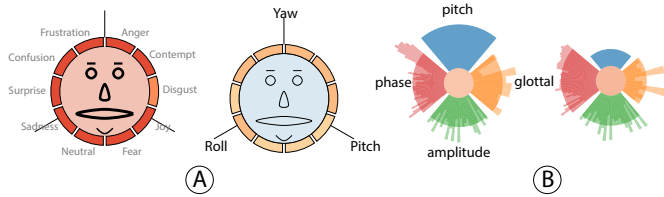
Fig. 4. The glyph designs in the *Projection View*. A: Chernoff face glyph designs. The left one with darked colored rings and thick strokes of face parts indicates intense facial movement, while the right one suggests little facial movement. B: Audio glyph designs. The left one with big blue sectors indicates high pitch, while the right one suggests low pitch.

distribution of importance of all three modalities. Their vertical orders show the total influences of the corresponding modalities, which are summed up by the blue bars to the left. The color of the bars inside the barcodes represents the importance values.

Besides, between two neighboring layers, links are drawn from the parent nodes to their child nodes. The width of a link is proportional to the importance of the child node to the model predictions.

**Design choices.** We have considered an alternative design (Fig. 3A) based on the Sankey diagram to reveal the intra- and inter-modal interactions and their importance to the predictions. It consists of three parts, the ground truth information at the left, the influences of individual modality at the center, and the inter-modal interactions at the right. The width of a flow is proportional to the importance of the target node of the flow. The barcode chart of each node further displays the importance distribution. In addition, the orange lines of the nodes show the error distribution to guide the exploration. However, one expert commended that it would be necessary to demonstrate more detailed information on each node. For example, what modalities dominate the predictions, and what is the frequency? Therefore, we augment the nodes with graphs and further convert the Sankey diagram into a compact tree-like layout, which leads to the current design (Fig. 3B).

### 5.3.2 Template View

To facilitate the exploration of feature sets and their influences, the *Template View* (Fig. 1C) summarizes frequent and influential templates of multimodal features in a table (*Q6*).

**Visual designs.** The *Template View* has four columns describing information about the template types, support, importance, and predictions and errors (*R1, Q6*). The first column records the names of feature sets by default. If a feature set contains frequent and important features, a green bar will be placed to the right denoting the number of children for the feature set. Users can collapse the corresponding row for detail by clicking the ⊞. The second table column displays the frequency for the templates. The distribution of the templates' importance and prediction information is visualized in the third and fourth columns. They share the same visual representations with the *Summary View* (Sect. 5.3.1). Users are enabled to sort the templates according to their support, importance, and errors. In such a way, they can prioritize their efforts in diagnosing the complex model behavior.

### 5.3.3 Projection View

To further support the subset-level exploration of model behavior (*R3, Q6*), the *Projection View* (Fig. 1D) connects the multimodal feature templates in the *Template View* with the instances. It allows users to examine the detailed information (e.g., feature values, prediction errors) about features across the instances. For example, after users select the "ADJ" template in the *Template View*, they may feel intrigued by what adjectives associate with large errors or with positive predictions. Then, they need to further inspect the individual instances.

**Visual designs.** To summarize the feature sets of a group of instances, we project the high-dimensional features onto a 2D plane using t-SNE [71]. Thus, instances with similar features will be placed close to each other. Given textual, acoustic, and visual features are heterogeneous, we design three different glyphs to encode the feature sets of the instances. Users can switch between views to see the feature distribution of each modality. Moreover, to help diagnose the model

behavior (e.g., errors), we add a heatmap as the background to display the distribution of prediction errors or template importance.

- *Language*: since words already carry semantic meanings, we use them to represent the textual features. In addition, we add a circle for each word, whose color encodes the sentiment prediction.
- *Vision*: our glyph designs for facial features (Fig. 4A) are inspired by Chernoff face [12], which is popular for displaying facial expressions. However, the original Chernoff face cannot reflect information such as head movement. Therefore, we add three sticks around the face to indicate the head movement in the yaw, pitch, and roll axis, respectively. The outer ring encodes the whole face information (e.g., emotion in our case), where the dark color suggests large feature values. Moreover, the stroke width of face parts (e.g., nose) and sticks mean movement intensity. The sentiment prediction is revealed by the face's background color.
- *Audio*: to help understand acoustic features, we group them into higher-level classes (Sect. 5.2.2). As shown in Fig. 4B, each colored sector represents the features of a class, where the radius relates to feature values. The sectors at the front summarized the average values of normalized features, while the small ones at the back display detailed feature values of the classes. Additionally, the inner circle color shows the sentiment prediction.

### 5.3.4 Instance View

The *Instance View* (Fig. 1E) provides local explanations by visualizing the important multimodal features and the context (i.e., transcripts and videos) of individual instances (*Q7*).

**Visual designs.** The left column presents a visual summary of the influences of modalities on the model predictions, as well as the prediction errors. Users can sort the instances according to different criteria (e.g., error) at the header and prioritize their efforts in instance-level exploration. In each row, the horizontal axes demonstrate the sentiment range, where the prediction and ground truth are marked. Between the two values, the thick red line suggests the error. Below the prediction mark, three colored rectangles represent the aggregated feature importance values of the three modalities. The length and color of each rectangle encode the magnitude and sign of importance. For example, the modality with negative influences on the prediction will be encoded by a blue rectangle and placed at the right. In addition, the feature table below allows users to sort and search for the importance values of features or modalities.

To promote a comprehensive understanding of the context of individual instances, the right column highlights the important features of the instances. Unlike intuitive texts, the acoustic and visual features are harder to recognize. Thus, we align them with the spoken words and draw the most important ones using orange lines. The lines above the words correspond to acoustic features, while the lines below represent the visual features. The vertical offset of the lines denotes the feature values, and hence the fluctuations indicate the feature variation. In addition, the backgrounds of texts or feature lines reflect the importance of multimodal features at a word level.

The *Instance View* also provides video context for instance-level exploration. When users click on the rows of the table, the corresponding video clips will pop up and play. To make the visual features more intuitive, the top-ranked facial features (sorted according to importance value) are highlighted with bounding boxes that cover the corresponding parts of the face. Users can further find the detailed facial action units and their concrete meanings by hovering on the boxes.

### 5.3.5 User Interactions

The *M²Lens* provides a rich set of interactions, which help unify the different views and facilitate multi-level and multi-faceted exploration with details on demand.

**Brushing.** Users can brush the barcodes in the last layer of *Summary View* to emit a query on the specific data instances of an interaction type. Then, the *Template View* and *Instance View* will show the related templates and local explanations, respectively.

**Clicking.** Many interactions in the system can be triggered and undone by clicking. For example, clicking the table rows in the *Tem-*

*plate View* will filter the irrelevant instances in the *Projection View* and *Instance View*. Users can switch between feature projections of different modalities by clicking the radio buttons in the *Projection View*. When clicking the table rows in the *Instance View*, the corresponding instances in the *Projection View* will be shown, and its video clips will pop up and play. In addition, users can click on the header of the *Template View* and *Projection View* to undo the previous selections.

**Lasso and semantic zooming.** To facilitate scalable exploration, users can use lasso or semantic zoom to focus on specific instances of interest in the *Projection View*. Then, the detailed information will be displayed in the *Instance View*.

**Searching, sorting, and filtering.** To narrow down the exploration space, users can sort and search for the instances or features in the table of *Template View* and *Instance View*. By adjusting the sliders in the *Projection View*, users can filter the instances according to the sentiment predictions and the feature importance of specific modalities.

## 6 EVALUATION

In this section, we demonstrate how $M^2Lens$ helps users understand and diagnose multimodal models for sentiment analysis through two case studies and interviews with three domain experts (*E1*, *E2*, and *E3*) using the CMU-MOSEI dataset. *E1* and *E2* are NLP researchers who have multiple top research publications on multimodal language analysis (e.g., emotion recognition). *E3* is a senior software engineer who has five years' experience in developing affective computing applications. The two cases are discovered by *E1* and *E2* during the system exploration in the interviews. The detailed feedback from all the experts is also collected and summarized.

### 6.1 Case One: Multimodal Transformer

In the first case, the expert *E1* explored and diagnosed a state-of-the-art model, Multimodal Transformer (MulT) [69], for sentiment analysis using the CMU-MOSEI dataset (Sect. 3.1). MulT fuses multimodal inputs with cross-modal transformers for all pairs of modalities, which learn the mappings between the source modality and target modality (e.g., vision → text). Then, the results are passed to sequence models (i.e., self-attention transformers) for final predictions. All the multimodal features of the input data are aligned at the word level based on the word timestamps. Following the settings of previous work [69], we trained, validated, and evaluated MulT with the same data splits (training: 16,265, validation: 1,869, and testing: 4,643). The details about the MulT are included in the supplementary material.

During the exploration, *E1* observed that the language modality often dominates the predictions, and the model cannot handle the negations in sentiment analysis very well. He further investigated the dominance of visual modality, where "Joy" and "Sadness" (two facial emotions) frequently co-occur. It was thought to be caused by the intense facial muscle movement, which was also captured by the model.

#### 6.1.1 Dominance of Language Modality

**Global summary (*R1, R2*)** After selecting the MulT and valid set in the *User Panel*, *E1* felt interested in how individual modalities and their interplay contribute to the model predictions. By looking at the second layer of the *Summary View* (Fig. 1B), *E1* found that the language modality (indicated by the letter "L") has the largest influence among the three modalities since it has the longest bar to the left and widest range of dots in the bee swarm plot. On the contrary, the acoustic modality (indicated by the letter "A"), which ranks at the bottom, has the least influence. Then, *E1* examined the last layer, where the *dominance* group with the widest barcode charts is shown at the top. Within the group, he discovered that the longest bars attach to the language modality, and the color of the prediction barcode aligns well with that of the language barcode. Thus, *E1* concluded that the language also plays a leading role in the *dominance* relationship. Furthermore, he noticed that there are a group of dense blue bars appearing at the end of the language barcode, where the errors are relatively large (as indicated by the yellow curve above the dashed line). He wondered what features or their combinations cause the high errors. Therefore, we brush the corresponding area of the blue bars.
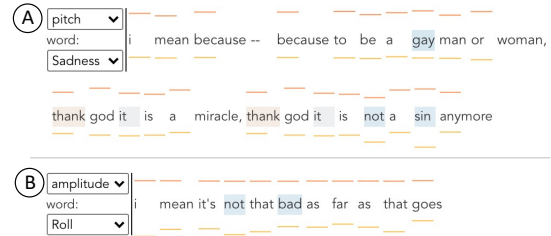


Fig. 5. Examples of double negations. "not...sin" (in A) and "not..bad" (in B) are considered as indicators for negative sentiment by the model. However, these phrases reflect sentiments that are slightly positive.

**Subset exploration (*R1, R3, R4*)** The *Template View* (Fig. 1C) lists all the frequent and important feature templates for the brushed instances in the *Summary View*. By sorting them in descending order of error, *E1* found that the "PRON + PART" appears at the top with one child feature. Then, he collapsed the row and found that 21 instances contain the word "not", where it negatively influences the predictions (blue dots in the bee swarm plot in the "importance" column). Next, he clicked "not" to see the details about this feature in the *Projection View*. Zooming in on the word "not", several similar negative words (e.g., "isn't", "wouldn't") were observed (Fig. 1D). They were all located in a red area, indicating large errors. *E1* speculated that the model cannot deal well with negations. Subsequently, he lassoed these words to closely examine the corresponding instances in the *Instance View*.

**Instance exploration (*R1, R3, R4*)** To further evaluate how the model handles negations, *E1* started with the instances with large errors in the table (Fig. 1E). When exploring the top-listed examples, *E1* observed that negations always have significant negative influences on the predictions, and the model fails to interpret the true sentiment. For example, *E1* found a case where the language modality dominates the negative sentiment prediction, and the word "not" is highlighted in blue (Fig. 1E). However, the true sentiment of this sentence is positive, where the starting phrase "I really like" demonstrates the positive attitude. However, the model fails to extract the keywords and relies on the negation (i.e., "not") to predict the negative sentiment. Moreover, *E1* noticed that when double negations appear in a sentence (Fig. 5), the model tends to treat them separately and regards both of them as indicators for negative sentiment. However, in fact, these double negations reflect sentiments that are slightly positive.

#### 6.1.2 Dominance of Visual Modality

**Global summary (*R1, R2*)** *E1* referred back to the "***dominance***" group in the *Summary View*, where a collection of red bars from the prediction barcode conform with the ones from the visual modality (highlighted red in Fig. 1B). The visual modality dominates the predictions, and the error line chart above suggests a low error rate in contrast with the previous case in Sect. 6.1.1. Motivated by this observation, *E1* brushed the red bars to investigate the patterns in the visual features.

**Subset exploration (*R1, R3, R4*)** In the *Template View*, "Face Emotion" has the largest support (Fig. 6A). After unfolding the row, *E1* found that "Joy + Sadness" is a frequent and important combination. This intrigued him to find out how a contrary emotion pair co-occurs. After clicking the template, the corresponding glyphs are highlighted in the *Projection View* (Fig. 6B). Most of them are found outside the red area, which verifies that the instances with "Joy + Sadness" often have small prediction errors. He decided to inspect these instances.

**Instance exploration (*R1, R3, R4*)** Through browsing the instances and their videos in the *Instance View*, "Joy" and "Sadness" are often considered important visual features with positive influences. Additionally, *E1* found their co-occurrences may be due to the presence of intense and rich facial expressions in the videos. These expressions generally involve the movement of the related facial action units in "Joy" and "Sadness". For example, after *E1* clicked on the instances, he noticed all the face parts (i.e., nose, eyes, brows, mouth, and chin) of the corresponding glyphs (Fig. 6B) in the *Projection View* has thick strokes, which suggests intense movements. When he watched the original videos, the bounding boxes of "Joy" and "Sadness" always popped up as important visual features. Hovering on the boxes and

Fig. 6. "Joy + Sadness" co-occurrence patterns. A: "Joy + Sadness" is a frequent and important feature template in the table. B: The raw video information and corresponding glyphs of three representative instances of the "Joy + Sadness" template.

examining the facial expressions and their explanations, *E1* concluded that the extreme facial expressions triggered the movement of the action units in "Joy" and "Sadness", and the model seemed to capture these important visual facial expressions.

Conclusions. During exploration, *E1* discovered that MulT cannot handle double negations very well, though it is a state-of-the-art model. He commented augmenting double negation examples or preprocessing them into positive forms can further improve the performance.

## 6.2 Case Two: EF-LSTM

In this case, the expert *E2* explored the popular RNN-based model, EF-LSTM [22], for multimodal sentiment analysis using the CMU-MOSEI dataset. The dataset setup and feature processing are the same as Case One (Sect. 6.1). EF-LSTM concatenates textual, acoustic, and visual features at each word. Then, it uses an LSTM model to derive the input representations for the predictions. The details of the model are provided in the supplementary material.

Through interactive explorations with M²Lens, *E2* was surprised to find that EF-LSTM does not learn sentiment in text. Also, he noticed that the acoustic modality has the largest influence on the sentiment prediction results among the three modalities, and the voice pitch always plays a negative role in the sentiment predictions.

### 6.2.1 No Meaningful Information Learned in Text

**Global summary (*R2*)** After selecting the valid set and EF-LSTM, *E2* started with the *Summary View* to gain an overview of the impacts of the modalities (Fig. 3B). By comparing the range of dots in the three bee swarm plots, he was surprised to find that acoustic modality is the most influential modality, then comes the language modality. In addition, the language modality always exhibits a positive impact on the sentiment. These findings are quite counter-intuitive. Thus, *E2* first explored text-related interactions by tracking the thickest links from the language modality to the third layer. He noticed that "*complement*" group shows at the top, and the text plays a leading role within the group. Then, he brushed the whole group to see textual feature patterns.

**Subset exploration (*R3, R4*)** The strange thing is that no textual templates and text glyphs were spotted in the *Template View* and the *Projection View*, respectively. *E2* suspected that the model does not learn any important language features (i.e., words) for sentiment analysis. Then, he referred to the *Instance View* to validate his doubt.



Fig. 7. Negative influences of voice pitch. A: "pitch" is the most frequent acoustic template, and it always has a negative impact (as indicated by the dots in the bee swarm plot). B: The selected group of instances with large pitch values (as indicated by the large radius of the blue sectors). C: Two high-error cases where the model captures the turning points of the pitch but wrongly associates pitch with negative influences.

| 1 | *It's run by a fantastic team of professors; they are always available for you.* |
| 2 | *(Umm) this movie was excellent.* |

**Instance exploration (*R1, R3, R4*)** When exploring the instances in the *Instance View*, *E2* found that the model fails to recognize potentially-important words for sentiment analysis, such as "fantastic" (in line #1), "excellent" (in line #2). None of them is highlighted with colors in the *Instance Detail*. *E2* also noticed every word of the sentences in the feature table has evenly low positive importance scores (less than 0.1). This explains why the language modality always has positive influences and further proves that the model does not capture the sentiment in text.

### 6.2.2 Negative Influences of Voice Pitch

**Global summary (*R1, R2*)** *E2* paid attention to the most influential modality (i.e., the acoustic modality) in the *Summary View* (Fig. 3B), where a negatively-skewed distribution of dots was shown. In addition, he noticed that within the "conflict" group, the acoustic modality plays a negative role (blue bars) throughout the time. Thus, *E2* brushed this group to investigate the negative influence of acoustic features.

**Subset exploration (*R1, R3, R4*)** *E2* found the "pitch" is the most frequent acoustic template in the *Template View* (Fig. 7A). Moreover, *E2* noticed that pitch always has a negative impact given the negatively-skewed distribution of dots in the third column. After clicking the row, he switched to the *Projection View* to see the pitch value distribution (Fig. 7B). He discovered that the acoustic glyphs are spread along a left-slanting line, where the radius of the blue sectors (i.e., pitch values) generally increases from left to right. Then, he selected a group of instances with the large pitch at the right corner for further inspection.

**Instance exploration (*R1, R3, R4*)** By browsing the instances and videos in the *Instance Summary* (Fig. 7C), *E2* observed that pitch is always the top important acoustic feature and is associated with negative influences. Although some important pitch variation signals in

the videos are captured by the model, he believed that the model is not reliable since it always regards the pitch as a strong negative sentiment indicator and he found many counterexamples. To name a few, in two cases (Fig. 7C), he found pitch ranks the first with negative importance in the feature table. And he noticed that some backgrounds of the orange lines (i.e., pitch values) are colored light blue (i.e., negative). By examining the offsets of all the orange lines, he thought the highlighted ones seem to be the turning points of pitch values. He speculated that the model captures the important signals in audio. He further checked the original video and verified the observations. However, the speakers sound high-spirited, and the pitch should reflect positive sentiment.

Conclusions. Through the case study, *E2* found that EF-LSTM seems not able to capture the sentiment in text. He reasoned that the simple early feature fusion may lead to textual information loss. He speculated that some more advanced model designs (e.g., transformer) can be incorporated into the model to facilitate text understanding. Given the negative impacts of voice pitch, *E2* thought that removing the pitch feature may increase the model accuracy.

### 6.3 Expert Interviews

We collected the feedback from the *one-on-one* interviews with the aforementioned three domain experts (*E1, E2, E3*). None of them have tried the system before the interviews. We first introduced the background and system designs. Then we asked the experts to use $M^2Lens$ to diagnose two state-of-the-art models (i.e., multimodal transformer and EF-LSTM) on the CMU-MOSEI dataset. After a 50-minute exploration, we collected their feedback about the system workflow, system designs, application scenarios, and improvement suggestions.

**System workflow.** All the experts confirmed the effectiveness of the system workflow of $M^2Lens$ in providing explanations for multimodal sentiment analysis models. They mentioned that they usually rely on performance metrics or instance-level feature importance measures for model evaluation, which does not provide many details and is unable to support an in-depth analysis. Our system supplements them with global- and subset-level explanations, which facilitates a comprehensive and systematic understanding of model behaviors. *E1* and *E3* praised that the interaction summaries (i.e., *dominance*, *complement*, and *conflict*) are impressive and very useful for revealing both the model behaviors and the multimodal data characteristics. *E3* mentioned if he finds some modalities are influential in predicting sentiment using $M^2Lens$, he can consider reducing the number of modalities without losing much performance when deploying the model to low-end devices. *E1* added that the feature templates help generalize the model error patterns. *E2* summarized that the system assisted him in discovering interesting insights into the models. For example, he was surprised that EF-LSTM seems to not capture any sentiment information from the text.

**Visual designs and interactions.** Overall, the experts confirmed that the visualizations are useful and still easy to understand, and interactions are smooth. The *Summary View* is most favored by the experts for a quick overview of the learned intra- and inter-modal interactions. The designs of *Projection View* are also appreciated by the experts. *E3* really liked the heatmap for showing the error and feature importance patterns. *E1* thought the face glyphs are very intuitive, and the interactions such as lasso and zoom are really helpful for the exploration of a large amount of data. Moreover, he valued the video playback and the realtime highlighting of face parts for raw video browsing. Nevertheless, *E1* and *E2* said that the *Instance View* is a little complex, visualizing lots of information. Additionally, the experts responded that it took them a while (about 20 minutes) to fully grasp all the components and functions in the system.

**Improvements.** The experts offered constructive suggestions for improvements. *E3* requested a bookmark function to save user interaction histories (e.g., selection of templates) for further review. *E1* suggested that the system can add a comparison module for exploring and comparing different models at the same time. During the exploration, *E2* and *E3* observed that some large model errors are caused by dataset errors (e.g., a mismatch between the video and transcript). They recommended that the system should support correcting dataset errors.

## 7 DISCUSSION

Here, we discuss $M^2Lens$ regarding generalizability, scalability, multi-level and multi-faceted exploratory analysis, and learning curve.

**Generalizability.** $M^2Lens$ was developed to visualize and explain multimodal models for sentiment analysis. We demonstrated our system through case studies on two state-of-the-art models using the CMU-MOSEI dataset. However, $M^2Lens$ can also be used to explain other multimodal models on different sentiment datasets based on the feature importance computed by post-hoc explainability techniques. Furthermore, the interaction types (i.e., dominance, complement, and conflict) and feature templates can summarize multimodal features from the global and subset levels in other multimodal language analyses. For example, for the multimodal emotion recognition task, the system can explain what are the dominant modalities when "angry" is predicted. The feature templates can summarize the frequent and influential feature sets for "angry" and facilitate the exploration of model behaviors.

**Scalability.** Our approach also has some scalability issues, which come from the automated algorithms and visual designs. The bottleneck of our computational cost is the feature attribution methods. We use SHAP to compute the feature importance. It took about 25 minutes to process 2,000 instances of the CMU-MOSEI validation set. To speed up the process, we can employ techniques such as feature clustering, data sampling, and parallel computing. For the visual designs, the visual clutter can occur in the *Projection View*, where multimodal instances are encoded with different glyphs. To reduce this issue, $M^2Lens$ enables filtering instances according to the feature importance and sentiment predictions. Moreover, users can use semantic zoom to focus on instances of interest, which alleviates the overlapping issues.

**Multi-level and multi-faceted exploratory analysis.** $M^2Lens$ provides multi-level and multi-faceted explanations on the behaviors of multimodal models for sentiment analysis. A general workflow for our target users (e.g., model users and researchers) starts with the *Summary View*, where the global summary of the influences of individual modalities and their interplay is displayed. Then, users can specify an interaction type. Its influential and frequent multimodal features will be summarized in the *Template View* and *Projection View*. Users can examine their error and importance patterns, which helps prioritize their efforts for the instance exploration in the *Instance View*.

**Learning curve.** According to the feedback from the expert interviews, the experts pointed out that it took them some time (usually a 20-min trial) before smoothly using our system since our system contains a few components. However, they said that $M^2Lens$ is very helpful for them to explore the models. Moreover, they have derived comprehensive insights into the model behaviors and are eager to use $M^2Lens$ for model understanding and diagnosis in the future.

## 8 CONCLUSION AND FUTURE WORK

In this paper, we presented $M^2Lens$, a visual analytics system to help users understand and diagnose multimodal models for sentiment analysis. $M^2Lens$ provides multi-level explanations on model behaviors from language, acoustic, and visual modalities. It features an augmented tree-like layout for a global understanding of learned intra- and inter-modal interactions. Moreover, the feature templates and visualization glyphs of multimodal features facilitate the exploration of a group of frequent and influential feature sets. Through two case studies and expert interviews, we demonstrated $M^2Lens$ can provide deep insights into the state-of-art multimodal models for sentiment analysis.

In the future, we plan to enhance our system usability by adding functions, such as model comparison, data error correction. Also, we would like to extend our system to other multimodal applications (e.g., emotion recognition). Further, more domain experts can be invited to further validate the usability and effectiveness of $M^2Lens$ with more datasets and models for sentiment analysis.

# REFERENCES

[1] Y. Ahn and Y.-R. Lin. Fairsight: Visual analytics for fairness in decision making. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):1086–1095, 2020.

[2] S. Amershi, M. Chickering, S. M. Drucker, B. Lee, P. Simard, and J. Suh. Modeltracker: Redesigning performance analysis tools for machine learning. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 337–346, 2015.

[3] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al. Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115, 2020.

[4] T. Baltrušaitis, C. Ahuja, and L.-P. Morency. Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):423–443, 2018.

[5] O. Bastani, C. Kim, and H. Bastani. Interpretability via model extraction. *arXiv preprint arXiv:1706.09773*, 2017.

[6] A. Bilal, A. Jourabloo, M. Ye, X. Liu, and L. Ren. Do convolutional neural networks learn class hierarchy? *IEEE Transactions on Visualization and Computer Graphics*, 24(1):152–162, 2018.

[7] M. Brooks, S. Amershi, B. Lee, S. M. Drucker, A. Kapoor, and P. Simard. Featureinsight: Visual support for error-driven feature ideation in text classification. In *IEEE Conference on Visual Analytics Science and Technology*, pp. 105–112, 2015.

[8] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4):335–359, 2008.

[9] Á. A. Cabrera, W. Epperson, F. Hohman, M. Kahng, J. Morgenstern, and D. H. Chau. Fairvis: Visual analytics for discovering intersectional bias in machine learning. In *IEEE Conference on Visual Analytics Science and Technology*, pp. 46–56, 2019.

[10] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832, 2019.

[11] M. Chen, S. Wang, P. P. Liang, T. Baltrušaitis, A. Zadeh, and L.-P. Morency. Multimodal sentiment analysis with word-level fusion and reinforcement learning. In *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, pp. 163–171, 2017.

[12] H. Chernoff. The use of faces to represent points in k-dimensional space graphically. *Journal of the American Statistical Association*, 68(342):361–368, 1973.

[13] G. Degottex, J. Kane, T. Drugman, T. Raitio, and S. Scherer. COVAREP - A collaborative voice analysis repository for speech technologies. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 960–964, 2014.

[14] R. Ekman. *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, 1997.

[15] E. Friesen and P. Ekman. *Facial action coding system: A technique for the measurement of facial movement*. Consulting Psychologists Press, 1978.

[16] S. Gehrmann, H. Strobelt, and A. M. Rush. GLTR: Statistical detection and visualization of generated text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 111–116, 2019.

[17] M. Glodek, S. Tschechne, G. Layher, M. Schels, T. Brosch, S. Scherer, M. Kächele, M. Schmidt, H. Neumann, G. Palm, et al. Multiple classifier systems for the classification of audio-visual emotional states. In *International Conference on Affective Computing and Intelligent Interaction*, pp. 359–368. Springer, 2011.

[18] B. Haasdonk. Feature space interpretation of svms with indefinite kernels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(4):482–492, 2005.

[19] J. Han, J. Pei, Y. Yin, and R. Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery*, 8(1):53–87, 2004.

[20] A. W. Harley. An interactive node-link visualization of convolutional neural networks. In *International Symposium on Visual Computing*, pp. 867–877. Springer, 2015.

[21] A. Henelius, K. Puolamäki, H. Boström, L. Asker, and P. Papapetrou. A peek into the black box: exploring classifiers by randomization. *Data Mining and Knowledge Discovery*, 28(5):1503–1529, 2014.

[22] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[23] F. Hohman, M. Kahng, R. Pienta, and D. H. Chau. Visual analytics in deep learning: An interrogative survey for the next frontiers. *IEEE Transactions on Visualization and Computer Graphics*, 25(8):2674–2693, 2018.

[24] A. Hu and S. Flaxman. Multimodal sentiment analysis to explore the structure of emotions. In *In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 350–358, 2018.

[25] U. Johansson, L. Niklasson, and R. König. Accuracy vs. comprehensibility in data mining models. *Information Fusion*, 1:295–300, 2004.

[26] M. Kahng, P. Y. Andrews, A. Kalro, and D. H. Chau. Activis: Visual exploration of industry-scale deep neural network models. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):88–97, 2017.

[27] A. Karpathy, J. Johnson, and L. Fei-Fei. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*, 2015.

[28] P. W. Koh and P. Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1885–1894. PMLR, 2017.

[29] R. Konig, U. Johansson, and L. Niklasson. G-REX: A versatile framework for evolutionary data mining. In *Workshops Proceedings of the 8th IEEE International Conference on Data Mining*, pp. 971–974, 2008.

[30] V. Krakovna and F. Doshi-Velez. Increasing the interpretability of recurrent neural networks using hidden markov models. *arXiv preprint arXiv:1606.05320*, 2016.

[31] J. Krause, A. Perer, and E. Bertini. Infuse: Interactive feature selection for predictive modeling of high dimensional data. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1614–1623, 2014.

[32] J. Krause, A. Perer, and K. Ng. Interacting with predictions: Visual inspection of black-box machine learning models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 5686–5697, 2016.

[33] J. Krause, A. Perer, and H. Stavropoulos. Supporting iterative cohort construction with visual temporal queries. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):91–100, 2015.

[34] T. Kulesza, M. Burnett, W.-K. Wong, and S. Stumpf. Principles of explanatory debugging to personalize interactive machine learning. In *Proceedings of the 20th International Conference on Intelligent User Interfaces*, pp. 126–137, 2015.

[35] A. Lazaridou, N. T. Pham, and M. Baroni. Combining language and vision with a multimodal skip-gram model. In *Proceedings of the North American Chapter of the Association for Computational Linguistics*, pp. 153–163, 2015.

[36] Z. Liu, Y. Shen, V. B. Lakshminarasimhan, P. P. Liang, A. Zadeh, and L.-P. Morency. Efficient low-rank multimodal fusion with modality-specific factors. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 2247–2256, 2018.

[37] L. Longo, R. Goebel, F. Lecue, P. Kieseberg, and A. Holzinger. Explainable artificial intelligence: Concepts, applications, research challenges and visions. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, pp. 1–16. Springer, 2020.

[38] S. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pp. 4765–4774, 2017.

[39] C. Molnar. *Interpretable Machine Learning*. 2019. https://christophm.github.io/interpretable-ml-book/.

[40] L.-P. Morency, R. Mihalcea, and P. Doshi. Towards multimodal sentiment analysis: Harvesting opinions from the web. In *Proceedings of the 13th International Conference on Multimodal Interfaces*, pp. 169–176. ACM, 2011.

[41] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng. Multimodal deep learning. In *Proceedings of the 28th International Conference on Machine Learning*, pp. 689–696. Omnipress, 2011.

[42] B. Nojavanasghari, D. Gopinath, J. Koushik, T. Baltrušaitis, and L.-P. Morency. Deep multimodal fusion for persuasiveness prediction. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pp. 284–288, 2016.

[43] C. Olah, A. Mordvintsev, and L. Schubert. Feature visualization. *Distill*, 2017.

[44] K. Patel, J. Fogarty, J. A. Landay, and B. Harrison. Investigating statistical machine learning as a tool for software development. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 667–676, 2008.

[45] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1532–1543, 2014.

[46] H. Pham, P. P. Liang, T. Manzini, L.-P. Morency, and B. Póczos. Found in translation: Learning robust joint representations by cyclic translations between modalities. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 6892–6899, 2019.

[47] S. Poria, E. Cambria, R. Bajpai, and A. Hussain. A review of affective computing: From unimodal analysis to multimodal fusion. *Information Fusion*, 37:98–125, 2017.

[48] S. Poria, E. Cambria, D. Hazarika, N. Majumder, A. Zadeh, and L.-P. Morency. Context-dependent sentiment analysis in user-generated videos. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pp. 873–883, 2017.

[49] G. Potamianos, C. Neti, G. Gravier, A. Garg, and A. W. Senior. Recent advances in the automatic recognition of audiovisual speech. *Proceedings of the IEEE*, 91(9):1306–1326, 2003.

[50] W. Rahman, M. Hasan, S. Lee, A. Zadeh, C. Mao, L.-P. Morency, E. Hoque, et al. Integrating multimodal information in large pretrained transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2359–2369, 2020.

[51] S. S. Rajagopalan, L.-P. Morency, T. Baltrusaitis, and R. Goecke. Extending long short-term memory for multi-view structured learning. In *European Conference on Computer Vision*, pp. 338–353. Springer, 2016.

[52] G. A. Ramirez, T. Baltrušaitis, and L.-P. Morency. Modeling latent discriminative dynamic of multi-dimensional affective signals. In *International Conference on Affective Computing and Intelligent Interaction*, pp. 396–406. Springer, 2011.

[53] D. Ren, S. Amershi, B. Lee, J. Suh, and J. D. Williams. Squares: Supporting interactive performance analysis for multiclass classifiers. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):61–70, 2016.

[54] M. T. Ribeiro, S. Singh, and C. Guestrin. "Why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144, 2016.

[55] M. T. Ribeiro, S. Singh, and C. Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, pp. 1527–1535, 2018.

[56] M. T. Ribeiro, T. Wu, C. Guestrin, and S. Singh. Beyond accuracy: Behavioral testing of nlp models with checklist. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4902–4912, 2020.

[57] M. Robnik-Šikonja and I. Kononenko. Explaining classifications for individual instances. *IEEE Transactions on Knowledge and Data Engineering*, 20(5):589–600, 2008.

[58] A. Rogers, O. Kovaleva, and A. Rumshisky. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866, 2020.

[59] S. Rubin, F. Berthouzoz, G. J. Mysore, W. Li, and M. Agrawala. Content-based tools for editing audio stories. In *In Proceedings of the 26th Annual ACM Symposium on User Interface Software and Technology*, pp. 113–122, 2013.

[60] E. Sariyanidi, H. Gunes, and A. Cavallaro. Automatic analysis of facial affect: A survey of registration, representation, and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(6):1113–1133, 2014.

[61] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 618–626, 2017.

[62] L. S. Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2(28):307–317, 1953.

[63] D. Smilkov, S. Carter, D. Sculley, F. B. Viégas, and M. Wattenberg. Direct-manipulation visualization of deep networks. *arXiv preprint arXiv:1708.03788*, 2017.

[64] M. Soleymani, D. Garcia, B. Jou, B. Schuller, S.-F. Chang, and M. Pantic. A survey of multimodal sentiment analysis. *Image and Vision Computing*, 65:3–14, 2017.

[65] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. In *International Conference on Learning Representations (Workshop Track)*, 2015.

[66] H. Strobelt, S. Gehrmann, H. Pfister, and A. M. Rush. Lstmvis: A tool for visual analysis of hidden state dynamics in recurrent neural networks.

*IEEE Transactions on Visualization and Computer Graphics*, 24(1):667–676, 2017.

[67] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.

[68] G. Tolomei, F. Silvestri, A. Haines, and M. Lalmas. Interpretable predictions of tree-based ensembles via actionable feature tweaking. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 465–474, 2017.

[69] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6558–6569, 2019.

[70] Y.-H. H. Tsai, M. Ma, M. Yang, R. Salakhutdinov, and L.-P. Morency. Multimodal routing: Improving local and global interpretability of multimodal language analysis. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1823–1833, 2020.

[71] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 11 2008.

[72] H. Wang, A. Meghawat, L.-P. Morency, and E. P. Xing. Select-additive learning: Improving generalization in multimodal sentiment analysis. In *IEEE International Conference on Multimedia and Expo*, pp. 949–954, 2017.

[73] X. Wang, H. Zeng, Y. Wang, A. Wu, Z. Sun, X. Ma, and H. Qu. Voicecoach: Interactive evidence-based training for voice modulation skills in public speaking. In *In Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pp. 1–12, 2020.

[74] J. Wexler, M. Pushkarna, T. Bolukbasi, M. Wattenberg, F. Viégas, and J. Wilson. The what-if tool: Interactive probing of machine learning models. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):56–65, 2019.

[75] J. Wiebe, T. Wilson, and C. Cardie. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 39(2):165–210, 2005.

[76] A. Zadeh, M. Chen, S. Poria, E. Cambria, and L.-P. Morency. Tensor fusion network for multimodal sentiment analysis. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1103–1114, 2017.

[77] A. Zadeh, P. P. Liang, N. Mazumder, S. Poria, E. Cambria, and L.-P. Morency. Memory fusion network for multi-view sequential learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, pp. 5634–5641, 2018.

[78] A. Zadeh, P. P. Liang, S. Poria, P. Vij, E. Cambria, and L.-P. Morency. Multi-attention recurrent network for human communication comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, pp. 5642–5649, 2018.

[79] A. B. Zadeh, P. P. Liang, S. Poria, E. Cambria, and L.-P. Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 2236–2246, 2018.

[80] H. Zeng, X. Shu, Y. Wang, Y. Wang, L. Zhang, T.-C. Pong, and H. Qu. Emotioncues: Emotion-oriented visual summarization of classroom videos. *IEEE Transactions on Visualization and Computer Graphics*, 27(7):3168–3181, 2020.

[81] H. Zeng, X. Wang, A. Wu, Y. Wang, Q. Li, A. Endert, and H. Qu. Emoco: Visual analysis of emotion coherence in presentation videos. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):927–937, 2019.

[82] Z. Zeng, M. Pantic, G. I. Roisman, and T. S. Huang. A survey of affect recognition methods: Audio, visual, and spontaneous expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(1):39–58, 2008.

[83] J. Zhang, Y. Wang, P. Molino, L. Li, and D. S. Ebert. Manifold: A model-agnostic framework for interpretation and diagnosis of machine learning models. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):364–373, 2018.