

DynAMoS: The Dynamic Affective Movie Clip Database for Subjectivity Analysis

Author Name 1
Department Name
Affiliation Name
Affiliation Location
Author Email or ORCID

Author Name 2
Department Name
Affiliation Name
Affiliation Location
Author Email or ORCID

Author Name 3
Department Name
Affiliation Name
Affiliation Location
Author Email or ORCID

Abstract—In this paper, we describe the design, collection, and validation of a new video database that includes holistic and dynamic emotion ratings from 83 participants watching 22 affective movie clips. In contrast to previous work in Affective Computing, which pursued a single “ground truth” label for the affective content of each moment of each video (e.g., by averaging the ratings of 2 to 7 trained participants), we embrace the subjectivity inherent to emotional perceptions and provide the full distribution of all participants’ ratings (with an average of 76.7 raters per video). We argue that this choice represents a paradigm shift with the potential to unlock new research directions, generate new hypotheses, and inspire novel methods in the Affective Computing community. We also describe several interdisciplinary use cases for the database: to provide dynamic norms for emotion elicitation studies (e.g., in psychology, medicine, and neuroscience), to train and test affective content analysis algorithms (e.g., for dynamic emotion recognition, video summarization, and movie recommendation), and to study subjectivity and ambiguity in emotional reactions (e.g., to identify moments of emotional ambiguity or ambivalence within movies, identify predictors of subjectivity, and develop personalized/subjective affective content analysis algorithms). The database is made freely available to researchers for noncommercial use at MASKED FOR REVIEW.

Index Terms—database, emotion elicitation, content analysis, affective computing, subjectivity, multimodal, movie clips

I. INTRODUCTION

In the affective and medical sciences, it is common for researchers to use standardized stimuli such as text vignettes, images, audio clips, and video clips to (try to) elicit/induce emotions in their participants [1]. This methodology can be used to study the effects of emotion on various psychological and physiological processes, to identify correlates of and group differences in emotional reactivity, and even to detect and quantify affective dysfunction in individual participants. The focus in such studies is on the affective experiences of the *participants* themselves (i.e., what they feel in response to the stimuli) as opposed to their perceptions of others (e.g., what the people portrayed in the stimuli seem to be feeling).

Movie clips are popular stimuli in such studies because movies are often expertly crafted (e.g., by actors, directors, sound designers, and editors) to produce a wide range of

emotional reactions in viewers. Due to their multimodal nature and extended duration, they tend to elicit stronger and more complex emotional responses than other stimulus types [2].

Participants are typically asked to provide a *holistic* report of their emotional reaction after watching each movie clip (e.g., in terms of *discrete* categories like anger and sadness or *continuous* dimensions like valence and arousal [3]). However, participants’ emotions often evolve over time during a movie clip. To address this issue, specialized methods have been developed to collect *dynamic* reports during stimulus presentation (e.g., by having the participant move a dial or lever to indicate changes in emotion) [4]. The resulting time-series (sometimes called ‘traces’) can capture the unfolding of emotional reactions over time with high granularity [5].

In the Affective Computing community, dynamic ratings of experienced and perceived emotion have been used to create ‘ground truth’ labels for emotion recognition and affective content analysis [3]. However, because each rater is a unique individual with their own history, background, and constellation of affective traits, there is inevitably some degree of inter-rater variability or ‘subjectivity’ in their ratings of each stimulus. This variability is often considered a nuisance—a source of noise to be minimized, e.g., by averaging across raters, training raters to consensus, or switching from continuous rating scales to ordinal rankings [6]. However, we contend that inter-rater subjectivity is actually a fascinating phenomenon worthy of academic study in its own right.

Embracing the existence of this subjectivity leads to many intriguing questions. Why do different participants experience the same stimuli in such different ways? How structured and predictable are individual participant’s responses? Are some stimuli (or parts of stimuli) associated with greater inter-participant variability than others, and if so, how well can we estimate their degree of subjectivity based on their content alone? We believe that the Affective Computing community is well-suited to begin answering these questions; however, doing so will require new datasets and novel methods [7].

In service of this goal, we present the DynAMoS database, which we designed to facilitate research on the dynamic and subjective aspects of emotional reactions to movie clips.

This work was partially supported by grant _____ from the _____. The views reported in this manuscript are the authors’ and do not necessarily reflect those of the funding agency.

A. Previous Databases

Given space constraints, we focus our literature review on video databases that provide *dynamic* emotion ratings. We use n to refer to the number of videos in a database and k to refer to the number of participants who rated each video.

The USC CreativeIT database [8], [9] includes videos of English-speaking theatrical performers ($n = 100$) improvising dyadic scenes ranging from 2 min to 10 min. Performers' expressed emotion (i.e., activation, valence, and dominance) was dynamically rated by observers ($\bar{k} = 3.2$).

The SEMAINE database [10] includes videos of English-speaking participants ($n = 24$) interacting with human-controlled virtual agents for 5 min. Participants' expressed emotion (i.e., arousal, expectation, power, and valence) was dynamically rated by observers ($\bar{k} = 6$) [11].

The AVID-Corpus [12] includes videos of German-speaking participants ($n = 292$) in directed speech tasks for 20 min to 50 min. Participants' expressed emotion (i.e., valence and arousal) was dynamically rated by one observer ($k = 1$).

The RECOLA database [13] includes videos of French-speaking participants ($n = 46$) in a computer-mediated remote collaboration task. Participants' expressed emotion (i.e., valence and arousal) during the first 5 min of the task was dynamically rated by observers ($k = 6$).

The DECAF database [14] includes a 51 s to 128 s clip from numerous English-language films ($n = 36$). Multiple participants ($k = 7$) dynamically rated each clip in terms of their own experienced emotion (i.e., valence and arousal).

The COGNIMUSE database [15] includes a 30 min video clip from several English-language films ($n = 7$). Each clip was dynamically rated by multiple participants ($k = 7$) in terms of the emotion (i.e., valence and arousal) they thought the director intended to convey, the emotion they themselves experienced in response to the clip, and the emotion they expected other viewers would experience.

The AFEW-VA database [16] includes short (0.4 s to 5.0 s) video clips ($n = 600$) from English-language films. The expressed emotion (i.e., valence and arousal) of each video frame was rated by two observers ($k = 2$).

Note that, based on their papers, the HUMAINE [17] and BINED [18] databases also seem to have included dynamic emotion ratings but appear to no longer be accessible.

In summary, there are currently seven video databases that include dynamic emotion ratings. However, only two include ratings of participants' subjectively experienced emotion, with all others focusing on observers' perceptions of the emotion being expressed by the individual in the video. Additionally, all databases used a small number of raters with minimal diversity (e.g., 1 to 7 college students), which limits the reliability of their average ratings and the representativeness of their sample of raters to the population of all possible raters.

We contend that the affective sciences need a new database of videos with dynamic emotion ratings that focuses on raters' subjectively experienced emotion, includes a large number of diverse raters viewing multiple videos, and has videos that are longer than several seconds in duration (so that they capture

more complex dynamic processes) but shorter than 10 min (so that showing multiple clips to the same participant is feasible).

B. The Current Study

The current study presents a new database that begins to meet this need and fill in the gaps of previous databases. It includes a set of 2 min to 7 min clips from English-language movies ($n = 22$), and each clip has been dynamically and continuously rated in terms of the subjectively experienced emotion (i.e., valence) of a large and diverse group of observers ($\bar{k} = 78.7$). This number of raters is over 10 times the number used in previous databases and opens up exciting new possibilities into the study of affective experiences, subjectivity/ambiguity, and personalized modeling.

The contributions of the current study are threefold:

- 1) We provide a new database of affective movie clips with dynamic and holistic ratings of experienced emotion from a large and diverse sample of participants and discuss new applications for this type of data.
- 2) We provide a website that documents the data overall and for each clip, and also discuss the benefits of creating such a website for newly released databases.
- 3) We demonstrate the use of modern rating validation techniques including Bayesian generalizability studies, inter-rater reliability analysis (i.e., ICCs) with incomplete data, and coefficient categorical omega for the internal consistency of ordinal scale scores.

II. METHODS

A. Movie Clip Selection

A set of 22 clips from high-definition movies of different genres (e.g., comedy, romance, drama, action) was identified using open libraries, such as YouTube and the set of films used in the Human Connectome Project [19]. Selection was based on the following criteria: (1) each clip must contain dialogue spoken by real-life actors (not animated or computer generated characters) in largely camera-facing orientations, (2) the set of clips must evoke a range of emotional reactions spanning positive and negative valence of various intensity levels, and (3) the set of clips must represent a diversity of movie actor demographics, including in gender, age, and race/ethnicity. The clips were 2 min to 7 min long to preserve their naturally evolving narrative structure and time-varying emotional content, while keeping them in a range that makes showing multiple clips to the same participant feasible.

B. Movie Clip Processing

Each movie clip was extracted from a Blu-Ray copy of its source movie to a separate MPEG-4 file with a frame width of 1920 px, a frame rate of 23.976 fps, and an audio sampling rate of 48 kHz. English-language subtitles were extracted and written to a separate text file (in SubRip Text format), including timestamps for each line of dialogue. Formatting syntax (e.g., HTML code to color a subtitle) and information for the deaf and hard-of-hearing (e.g., subtitles that describe the background music) were manually removed from the 'raw' version of each subtitle file to create a 'transcript' version.

C. Participant Recruitment

Participants were recruited from the local community using the platform; they were thus all living in the USA at the time of participation. Applicants were first screened to confirm their eligibility for the study; inclusion criteria included (1) having no uncorrected sensory, cognitive, or emotional impairments, (2) being age 18–60 years old, (3) being fluent in English, and (4) having access to a laptop or desktop computer and a quiet environment for the study sessions. After signing a consent form, participants were asked about their demographic background (i.e., sex/gender, race, ethnicity, and age).

D. Participant Procedure

All testing procedures were conducted remotely via Zoom video conferencing. Participants completed one or two testing sessions, each 90 min long, in each of which they viewed and rated 11 movie clips presented in a randomized order. For each movie clip, the following three steps occurred. First, a brief and standardized description of the clip was read aloud by the experimenter; this description was meant to orient participants to the scene and provide any necessary contextual information. All descriptions are provided on the database website. Second, the participant watched the clip, without subtitles, while simultaneously providing continuous valence ratings (described below). Third, after the clip ended, the participant provided holistic emotion ratings about the clip overall and answered several other questions (e.g., about whether the audio and video playback worked properly and how familiar they were with the clip’s source movie). After each session, participants were compensated 25 USD for their effort; thus, participants could earn up to 50 USD in total. This compensation rate is about a dollar higher than the minimum hourly wage in the state the data was collected in.

E. Affect Rating Procedure

While watching each movie clip, participants provided *dynamic* valence ratings using the CARMA software [20]. As depicted in Figure 1, CARMA displayed the video next to a vertical rating scale (represented by a color gradient) that ranged from -4 (*very negative*) to $+4$ (*very positive*). Participants were instructed to move a slider up and down within this rating scale (using the computer mouse or arrow keys) while watching the clip to reflect how negative/unpleasant or positive/pleasant it made them feel from moment to moment. The slider was positioned at zero at the onset of the clip, and participants could adjust the slider at any time; CARMA queried the relative location of the slider at 30 Hz and then averaged all queried values within 1 s temporal bins. Our rationale for this aggregation step is that it smooths the time series, removing unintentional motion artifacts, and yields scores that better align with the response time needed to process each moment of the clip (e.g., sensorially, emotionally, and cognitively) and make an appropriate motor response.

After watching each clip, participants provided *holistic* affect ratings using the Short Positive and Negative Affect



Fig. 1. Screenshot of Continuous Rating Collection in CARMA [20]

Scale (S-PANAS) [21], [22]. Participants were instructed to think about their overall emotional reaction to the movie clip and rate it on five positive affect items (alert, determined, enthusiastic, excited, inspired) and five negative affect items (afraid, distressed, nervous, scared, upset) using an ordinal scale from 0 (*very slightly or not at all*) to 4 (*extremely*). Scores on these items were combined through averaging to yield scale scores for positive affect and negative affect.

F. Validation Procedure

We first excluded the ratings of any participants who reported that a clip’s audio or video did not playback properly. Given our interest in rater subjectivity and ambiguity, we did not exclude participants for being outliers in terms of their ratings, although some database users may choose to do so.

We then estimated inter-rater reliability of the valence ratings within each movie clip (after excluding the first 10 s of each clip for reasons described in §III-D) and of the holistic ratings across all movie clips. Specifically, two-way intraclass correlation coefficients (ICCs) were estimated from Bayesian generalizability studies [23] using the `varde` R package [24]. There are many formulations of the two-way ICC, but the most relevant here is the average-measures consistency ICC for incomplete data or $ICC(Q, \hat{k})$, which quantifies the reliability of the average of all available raters’ ratings [23]. We followed common heuristics [25] in considering ICC values between .50 and .75 to be “moderate,” values between .75 and .90 to be “good,” and values above .90 to be “excellent.”

From these same generalizability studies, we also calculated the percentage of rating variance that was accounted for by the differences between rater intercepts; this percentage can be considered a rough index of how much rater-to-rater subjectivity was present in the ratings for each clip.

We also estimated the inter-item reliability (or internal consistency) of the holistic scales for each movie clip. To do so, we estimated coefficient categorical omega or ω_{u-cat} [26], [27] for the Positive Affect and Negative Affect scales from ordinal confirmatory factor analysis models of each clip’s holistic ratings (using the `lavaan` [28] and `semTools`

[29] R packages). We followed common heuristics [30] in considering omega values above .75 to be acceptable.

G. Website Generation

As a form of rich documentation for the database, as well as a supplement to this manuscript, we used the Quarto technical and scientific publishing system¹ to create a website for the database within R. This website reads in the database files and generates summary statistics, tables, and figures for the database as a whole and for each movie clip individually. It also includes screenshots from movie clips, word clouds of the subtitles, and visualizations of the ratings. These pages are all parameterized reports, which means they can be quickly and easily updated as the database grows and changes. The website is hosted using GitHub Pages.²

III. RESULTS

A. Movie Clip Summary

The database includes 22 movie clips, each drawn from a different English-language movie. As shown in Table I, the source movies were released between the years of 1991 and 2018 and the clips ranged from 130s to 425s in duration ($M=251.1$, $SD=90.8$). Example video frames from the movie clips are shown in Figure 2; even from this small sample of images, it is possible to see the diversity of characters, settings, and lighting conditions represented in the database.

B. Participant Summary

We recruited 83 participants. In terms of sex/gender, 56 reported being Female (68 %) and 26 reported being Male (31 %). In terms of race, 43 reported being White (52 %), 22 reported being Asian (27 %), 12 reported being Black (15 %), and 5 reported being another race (6 %). In terms of ethnicity, 70 reported being non-Hispanic/Latino (84 %) and 11 reported being Hispanic/Latino (13 %). In terms of age, participants ranged from 18–59 years old ($M=28.8$, $SD=9.9$).

C. Validation Results

Of the 83 recruited participants, 77 (93 %) completed both sessions and the remaining 6 only completed the first session (before dropping out). Out of the 1826 possible participant-clip combinations, data from 1702 (93 %) were collected without issue, data from 93 (5 %) were not collected (due to participant dropout and clips being added partway through recruitment), and data from 31 (2 %) were excluded due to participants reporting issues with audio and/or video playback.

Estimates of the inter-rater reliability of the dynamic valence ratings for each movie clip are presented in Table I. The reliability of the average of all available ratings per bin was “excellent” (i.e., $ICC(Q, \hat{k}) > 0.90$) across all clips. (Detailed results with variance and interval estimates per clip are provided on the database website.) These results imply that the average rating per bin can be used with high confidence in all clips to represent the rating of a “typical”

participant. However, excellent average-measures inter-rater reliability does not imply that there was little subjectivity in the ratings of individual participants. To the contrary, rater intercepts explained a substantial amount of variance in the ratings, ranging from 21 % to 65 % (also shown in Table I), which implies considerable subjectivity in the continuous ratings. These two statements may at first seem to be contradictory but are in fact expected, as the former is largely due to the well-known benefits of averaging multiple ratings (which offsets the idiosyncratic aspects of individual ratings) [31], [32].

Estimates of the inter-rater reliability of the holistic ratings cannot be calculated per movie clip (since they are only provided once per clip) and instead must be calculated across all clips. The reliability of the average of all available ratings was “excellent” (i.e., $ICC(Q, \hat{k}) > 0.90$) for all items as well as for the Positive Affect scale ($ICC=.990$) and the Negative Affect scale ($ICC=.988$) scores. (Item-level results and interval estimates are provided on the database website.) These results imply that the average rating per clip can be used with high confidence to represent the rating of a “typical” participant. The amount of variance in the holistic ratings explained by rater intercepts was 24.5 % for the Positive Affect scale and 19.3 % for the Negative Affect scale. These results imply that there was relatively less subjectivity in the holistic ratings as compared to the dynamic valence ratings.

Estimates of the inter-item reliability of the holistic ratings (within raters) can be calculated per movie clip (since each rater provides five ratings per scale). The categorical omega estimates per movie clip ranged from .66 to .89 ($M=.81$, $SD=.06$) for the Positive Affect scale and from .61 to .92 ($M=.82$, $SD=.09$) for the Negative Affect scale. Thus, most (but not all) clips had acceptable internal consistency when using the holistic ratings of any single participant. However, as described in the previous paragraph, the inter-rater reliability of the average across all available raters was excellent.

D. Dynamic Rating Summary

The distribution of dynamic valence ratings across all raters, bins, and movie clips is depicted in Figure 3. This distribution has a roughly Gaussian shape with ratings close to 0 being the most common and increasingly extreme ratings (in either direction) being increasingly less common.

Figure 4 displays the time series of dynamic valence ratings from four example movie clips (similar plots for all movie clips are available on the database website). The thick black line on each plot depicts the average of all available ratings per temporal bin; this is the highly reliable time series that users can use to represent the rating of a “typical” participant. The colored ribbons around the black line depict successively larger percentages of the ratings per bin, i.e., the yellow ribbon contains the most common ratings, the green bands show less common ratings, and the purple band shows even less common ratings. This novel visualization approach, which we call a “chromodoris plot” (after the colorful sea slugs of similar appearance), allows us to quickly see the central tendency of

¹<https://www.quarto.org>

²<https://pages.github.com>

TABLE I
SUMMARY INFORMATION ABOUT THE MOVIE CLIPS, SOURCE MOVIES, HOLISTIC RATINGS, AND CONTINUOUS RATINGS

Clip	Duration	Raters	Source Movie Information		Holistic PA		Holistic NA		Dynamic Valence		
			Title	Year	Mean	Omega	Mean	Omega	Mean	ICC	Rater
01	164s	81	Akeelah and the Bee	2006	2.49	.84	0.35	.78	1.51	.995	23.3%
02	162s	76	If Beale Street Could Talk	2018	1.19	.84	0.66	.86	0.57	.950	39.3%
03	274s	79	Catch Me If You Can	2002	1.16	.76	0.45	.85	0.30	.982	22.4%
04	132s	80	500 Days of Summer	2009	0.81	.78	0.42	.86	0.21	.947	50.7%
05	134s	79	Fences	2016	0.78	.72	1.79	.91	-2.06	.963	52.5%
06	218s	78	Forrest Gump	1994	0.96	.74	0.59	.78	-0.49	.932	47.9%
07	150s	76	Good Will Hunting	1997	0.92	.81	0.93	.88	0.16	.981	41.7%
08	239s	75	The Green Mile	1999	0.73	.78	2.24	.89	-2.57	.951	45.5%
09	425s	77	The King's Speech	2010	1.69	.88	0.80	.86	0.45	.946	49.0%
10	232s	73	Lady Bird	2017	0.99	.87	0.36	.71	0.18	.974	25.2%
11	414s	75	Legally Blonde	2001	2.20	.83	0.24	.72	1.15	.973	41.0%
12	130s	75	Little Miss Sunshine	2006	0.70	.66	0.46	.74	-0.45	.910	70.4%
13	240s	76	Marriage Story	2019	0.73	.88	1.49	.84	-2.18	.965	58.8%
14	341s	78	Miracle	2004	2.19	.89	0.31	.72	1.10	.994	21.1%
15	234s	75	Moonlight	2016	1.59	.86	0.47	.88	1.19	.969	46.2%
16	262s	72	No Country for Old Men	2007	0.83	.84	1.23	.92	-0.99	.954	58.2%
17	425s	78	The Parent Trap	1998	1.49	.84	0.19	.62	1.90	.968	38.7%
18	229s	79	Pulp Fiction	1994	0.95	.85	0.81	.91	0.21	.943	59.5%
19	237s	79	The Pursuit of Happyness	2006	2.01	.87	0.33	.68	1.35	.987	32.9%
20	274s	79	The Silence of the Lambs	1991	0.84	.72	1.57	.90	-1.37	.967	47.8%
21	294s	73	The Social Network	2010	0.91	.81	0.54	.79	-0.50	.923	40.7%
22	315s	75	Zodiac	2007	0.97	.79	1.85	.91	-0.99	.992	29.5%

Note. PA and NA = Positive and Negative Affect, Omega = Internal Consistency, ICC = Inter-Rater Reliability, Rater = Rater Variance.



Fig. 2. One Example Video Frame from Each Movie Clip (arranged left-to-right in the same order presented in Table I)

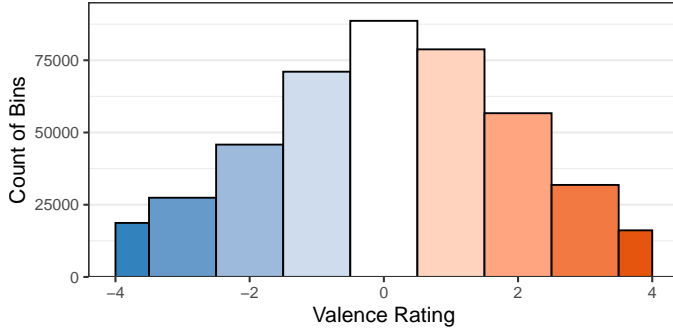


Fig. 3. Histogram of All Dynamic Valence Ratings

ratings as well as their spread (e.g., to locate parts of each movie clip that were more or less subjective).

A few insights can be derived from these visualizations. First, the ratings of the first 5s to 10s of each movie clip typically centered around zero with minimal spread. This pattern is likely due to the raters getting oriented to each clip. For many applications of the database (e.g., predicting ratings' mean or spread), it would make sense to exclude these bins from analysis. Second, some movie clips (e.g., *Akeelah and the Bee*) were relatively stable (or “stationary”) in terms of their

spread in ratings, whereas other clips had moments of sharp deviation. A striking example of the latter case is from *The Green Mile*; the ratings are quite negative throughout this clip, but at two moments (i.e., 01:00 to 01:30 and 3:00 to 3:40) the mean ratings became less negative and the spread in ratings increased dramatically. Similar increases in rating variability occur in both the *Fences* and *Lady Bird* time series, albeit to a lesser degree. The *Lady Bird* time series is especially interesting because large portions of the raters disagreed at many points whether the clip was positive or negative. These patterns raise several questions: What is it about these specific moments that leads to increased variability? Can we predict such moments from their content? What is it about the raters that caused them to respond differently to these moments? Can we predict such deviations from the norm?

E. Holistic Rating Summary

The distribution of holistic affect ratings for each movie clip is depicted in Figure 5. Clip averages ranged from 0.70 to 2.49 for Positive Affect ($M=1.23$, $SD=0.56$) and from 0.19 to 2.25 for Negative Affect ($M=0.82$, $SD=0.60$). Four clips exceeded 2.0 and nine clips exceeded 1.0 for average Positive Affect. In contrast, only one clip exceeded 2.0 and only six clips exceeded 1.0 for average Negative Affect.

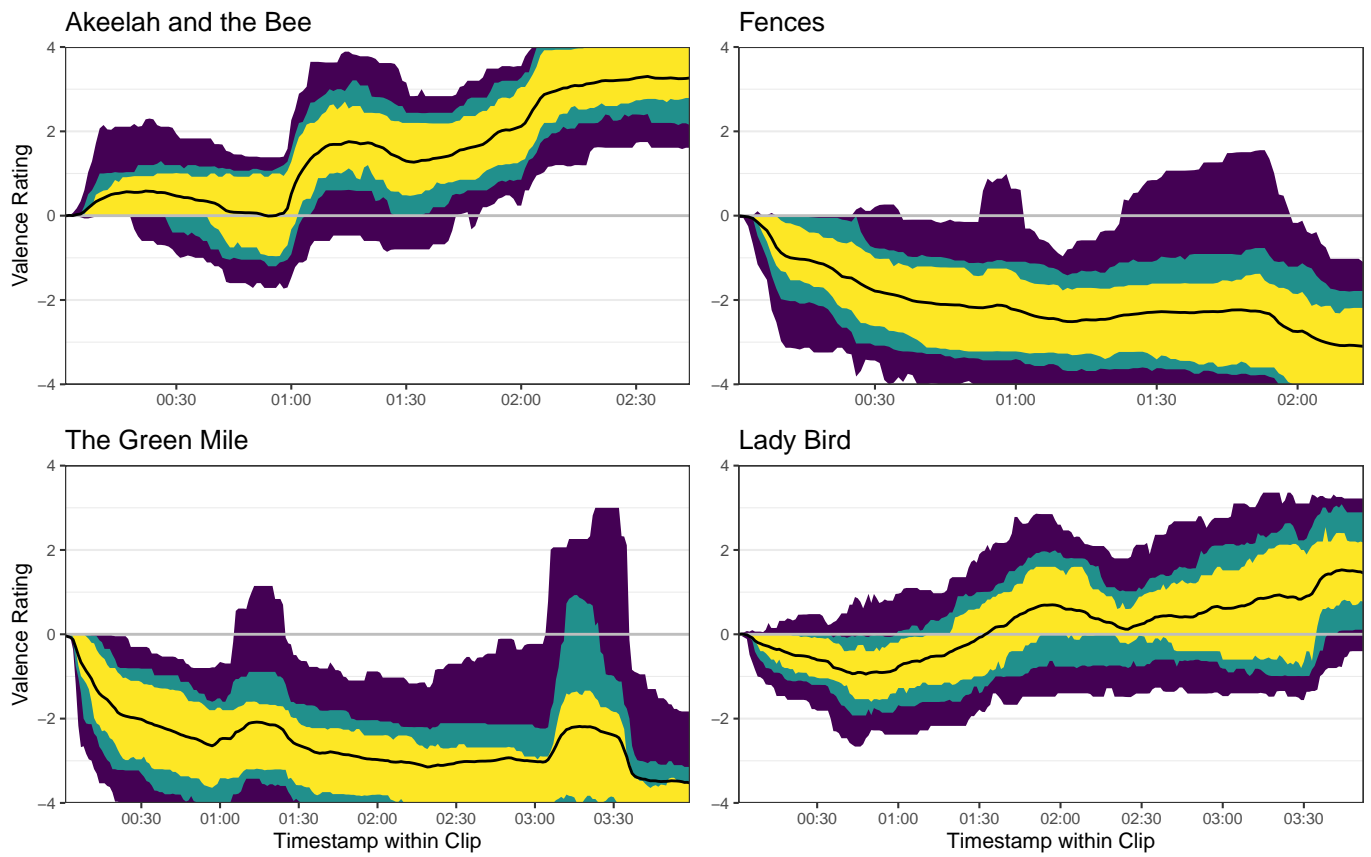


Fig. 4. Example Chromodoris Plots of Dynamic Ratings (Black = Mean Rating, Yellow = Inner 50%, Green = Inner 70%, Purple = Inner 90%)

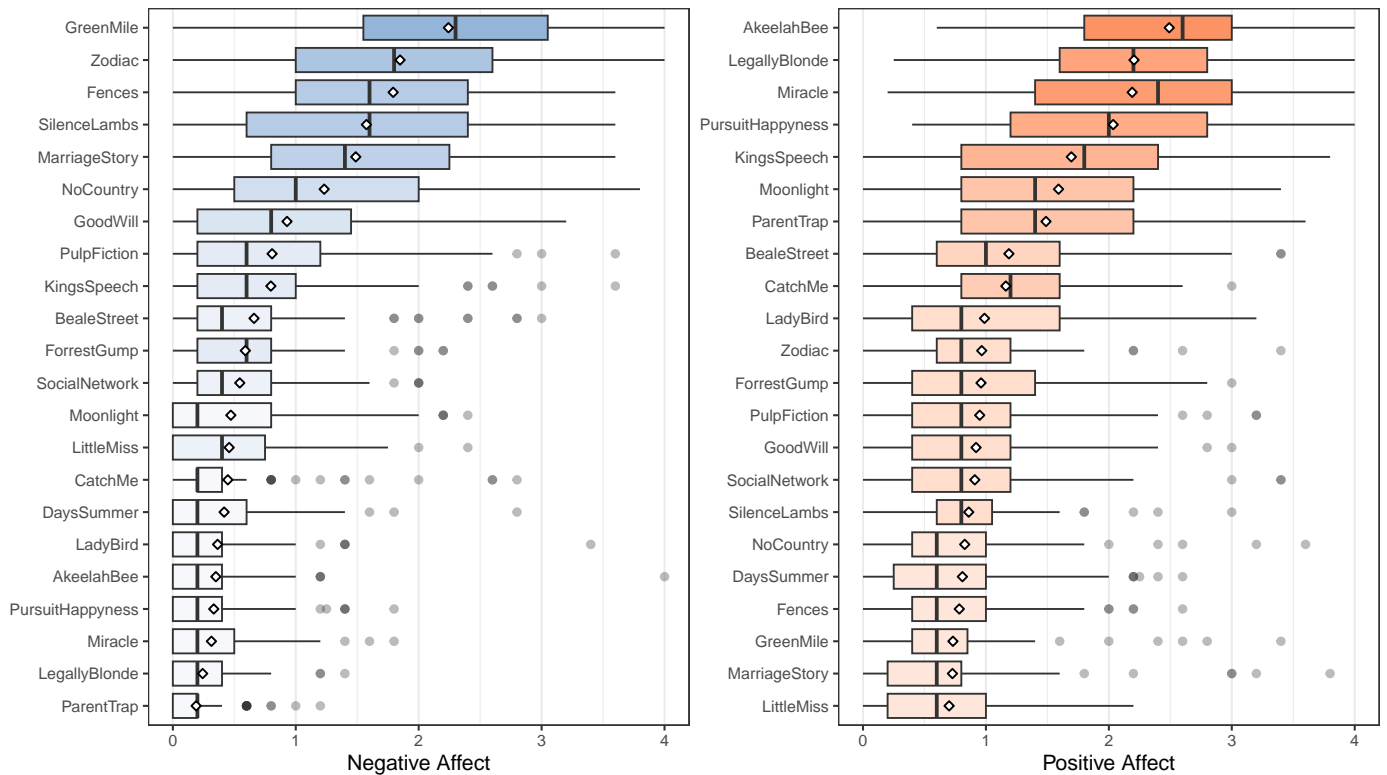


Fig. 5. Boxplots of the Distribution of Holistic Affect Ratings for Each Movie Clip (white diamonds = means, grey circles = outliers)

F. Website Summary

We were inspired by the websites of the SEMAINE database [10], which allows users to interactively filter and download different subsets of the data, and the GFT database [33], which allows users to access metadata, baseline results, and code. In creating the website for this database,³ we found it helpful to also include statistical summaries and visualizations of the data, thumbnails of frames from the clips, and scrollable tables with the subtitle lines. We found Quarto and GitHub Pages to be a powerful combination for this purpose and encourage other researchers to consider building similar documentation websites for future databases using these (or similar) tools.

IV. DISCUSSION

We argue that subjectivity and ambiguity are inherent to emotion representation and encourage researchers in Affective Computing to study these phenomena rather than simply trying to control for them. We argue that doing so has the potential to unlock new research directions, generate new hypotheses, and inspire novel methods. To promote work on this topic, we created and are sharing a new database with dynamic and holistic ratings of dozens of participants' emotional reactions to movie clips. In this paper, we describe the design, collection, and validation of the database. Results from our validation analyses support the trustworthiness of the database and reveal a large degree of subjectivity to be analyzed. We hope that this database will prove useful to researchers in several disciplines and will inspire more work on subjectivity and ambiguity.

A. Database Uses

The DynAMoS database has many potential use cases across several disciplines. First, it can be used as a standardized set of videos for emotion elicitation with normative data on emotional reactions (both holistic and dynamic). For example, these movie clips could be shown to new participants in psychology, medical, and neuroscience studies to induce positive and negative affect; furthermore, each new participant's ratings could be compared to the distribution of ratings in the database to quantify deviations from the norm.

Second, it can be used to train and test affective content analysis algorithms using traditional methods. For example, the multimodal information contained in each movie clip (e.g., images, speech, music, and subtitles) could be used to predict the holistic and/or dynamic ratings *averaged across raters*. Such predictions could be helpful for video summarization, movie recommendation, and identifying moments-of-interest.

Third, it can be used to study subjectivity and ambiguity [7] in affective experiences. For example, statistical or machine learning models could be used to predict the distributions of the holistic and/or dynamic ratings across raters. Possible predictor variables could include features of the movie clip (as in the second use case) and/or features of the participants themselves (e.g., demographics). Relatedly, this dataset may also be used in experiments on personalized/idiographic modeling [34] (e.g., predicting the ratings of specific individuals).

³Link to the database website is masked during peer-review.

B. Database Access

Summary information about the database is available on the database website³ and access to its full contents will be granted free-of-charge to researchers for noncommercial use. Potential users will need to request access through a form on the database website and sign a licensing agreement that states that they will (1) not use the data for commercial purposes such as developing products or patents, (2) not redistribute the data to unlicensed third parties, (3) not attempt to re-identify any participants who provided rating data, and (4) cite this paper in any published work making use of the database.

C. Limitations and Future Directions

Limitations of the current study include: (1) a focus on English-speaking movies and participants living in the USA, which limits generalizability to other languages and populations, (2) a relatively small number of movie clips, which limits how varied our set can be, (3) a focus on valence and the positive/negative activation model [35], which does not exhaustively capture the affective domain, and (4) relatively little information was collected about each participant, which limits our ability to study the sources of individual differences.

To address these limitations in future work, we (1) invite collaborations with researchers from other countries, (2) plan to add more movie clips that cover additional combinations of emotional content, actor demographics, spoken languages, and recording conditions, (3) plan to collect holistic ratings of discrete emotions and appraisal dimensions, and (4) plan to collect additional self-report measures of relevant characteristics such as personality and mental health.

We could also collect dynamic ratings of additional affective dimensions (e.g., arousal), but the costs of doing so would be non-trivial as it would require participants to either repeat the rating procedure or simultaneously rate multiple dimensions (which is possible but challenging [36], [37]). Also, at least in the case of arousal, separate rating may not be fully necessary as prior research suggests that arousal tends to increase with the intensity of both positive and negative emotion (i.e., with the magnitude or absolute value of valence) [35], [38].

We plan to release an extended version of this paper and database that adds multimodal features extracted from the movie clips, standardized partitions for cross-validation, and baseline predictive models for the use cases described in §IV-A. To promote interest in this area, we are also considering organizing workshops on subjectivity in Affective Computing.

ACKNOWLEDGEMENTS

We thank _____ for her support of this study, as well as _____ and _____ for their help in coordinating and testing study participants.

ETHICAL IMPACT STATEMENT

This work was approved by the governing Institutional Review Board prior to the work being carried out. Participants provided informed consent to complete the study and to have their deidentified data shared with other researchers. To further

protect the security and confidentiality of the data, we require users of the database to sign a licensing agreement (§IV-B).

We believe that the risk of our work having negative societal impacts is low, especially if its limitations (§IV-C) are appreciated by readers and users. The main ethical question we wrestled with regarded our use of clips from copyrighted movies. We believe that our use of brief clips to induce emotions during noncommercial research constitutes “fair use” according to §107 of the U.S. Copyright Act and is highly unlikely to harm the market for the copyrighted works.

REFERENCES

- [1] J. A. Coan and J. J. B. Allen, Eds., *Handbook of emotion elicitation and assessment*. Oxford University Press, 2007.
- [2] R. Westermann, K. Spies, G. Stahl, and F. W. Hesse, “Relative effectiveness and validity of mood induction procedures: a meta-analysis,” *European Journal of Social Psychology*, vol. 26, no. 4, pp. 557–580, 1996.
- [3] H. Gunes and B. W. Schuller, “Categorical and dimensional affect analysis in continuous input: Current trends and future directions,” *Image and Vision Computing*, vol. 31, no. 2, pp. 120–136, 2013.
- [4] A. M. Ruef and R. W. Levenson, “Continuous measurement of emotion: The affect rating dial,” in *Handbook of emotion elicitation and assessment*, J. A. Coan and J. J. B. Allen, Eds. Oxford University Press, 2007, pp. 286–297.
- [5] R. Cowie, E. Douglas-Cowie, S. Savvidou, E. McMahon, M. Sawey, and M. Schröder, “FEELTRACE: An instrument for recording perceived emotion in real time,” *ISCA Tutorial and Research Workshop on Speech and Emotion*, pp. 19–24, 2000.
- [6] H. Martinez, G. Yannakakis, and J. Hallam, “Don’t classify ratings of affect; rank them!” *IEEE Transactions on Affective Computing*, vol. 3045, no. c, pp. 1–1, 2014.
- [7] V. Sethu, E. M. Provost, J. Epps, C. Busso, N. Cummins, and S. Narayanan, “The Ambiguous World of Emotion Representation,” Sep. 2019, arXiv:1909.00360 [cs].
- [8] A. Metallinou, C.-C. Lee, C. Busso, S. Carnicke, and S. Narayanan, “The USC CreativeIT Database: A multimodal database of theatrical improvisation,” in *Proceedings of the 7th International Conference on Language Resources and Evaluation Workshops*, 2010, pp. 55–58.
- [9] A. Metallinou and S. Narayanan, “Annotation and processing of continuous emotional attributes: Challenges and opportunities,” in *Proceedings of the 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, Apr. 2013, pp. 1–8.
- [10] G. McKeown, M. Valstar, R. Cowie, M. Pantic, and M. Schröder, “The SEMAINE database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent,” *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 5–17, 2012.
- [11] B. Schuller, M. Valstar, R. Cowie, and M. Pantic, “AVEC 2012: the continuous audio/visual emotion challenge - an introduction,” in *Proceedings of the 14th ACM International Conference on Multimodal Interaction*, New York, NY, USA, 2012, pp. 361–362.
- [12] M. F. Valstar, B. W. Schuller, K. Smith, F. Eyben, B. Jiang, S. Bilakhia, S. Schneider, R. Cowie, and M. Pantic, “AVEC 2013: the continuous audio/visual emotion and depression recognition challenge,” in *Proceedings of the 3rd ACM international workshop on Audiovisual emotion challenge*, 2013, pp. 3–10.
- [13] F. Ringeval, A. Sonderegger, J. Sauer, and D. Lalanne, “Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions,” in *Proceedings of the 10th IEEE International Conference on Automatic Face and Gesture Recognition*, 2013, pp. 1–8.
- [14] M. K. Abadi, R. Subramanian, S. M. Kia, P. Avesani, I. Patras, and N. Sebe, “DECAF: MEG-Based Multimodal Database for Decoding Affective Physiological Responses,” *IEEE Transactions on Affective Computing*, vol. 6, no. 3, pp. 209–222, Jul. 2015.
- [15] A. Zlatintsi, P. Koutras, G. Evangelopoulos, N. Malandrakis, N. Efthymiou, K. Pastra, A. Potamianos, and P. Maragos, “COGN-IMUSE: a multimodal video database annotated with saliency, events, semantics and emotion with application to summarization,” *EURASIP Journal on Image and Video Processing*, vol. 2017, no. 1, p. 54, 2017.
- [16] J. Kossaifi, G. Tzimiropoulos, S. Todorovic, and M. Pantic, “AFEWA database for valence and arousal estimation in-the-wild,” *Image and Vision Computing*, vol. 65, pp. 23–36, 2017.
- [17] E. Douglas-Cowie, R. Cowie, I. Sneddon, C. Cox, O. Lowry, M. McRorie, J.-C. Martin, L. Devillers, S. Abrilian, A. Batliner, N. Amir, and K. Karpouzis, “The HUMAINE Database: Addressing the Collection and Annotation of Naturalistic and Induced Emotional Data,” in *Proceedings of the 2nd International Conference on Affective Computing and Intelligent Interaction*, 2007, pp. 488–500.
- [18] I. Sneddon, M. McRorie, G. McKeown, and J. Hanratty, “The Belfast induced natural emotion database,” *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 32–41, 2012.
- [19] D. C. Van Essen, S. M. Smith, D. M. Barch, T. E. J. Behrens, E. Yacoub, and K. Ugurbil, “The WU-Minn Human Connectome Project: An overview,” *NeuroImage*, vol. 80, pp. 62–79, 2013.
- [20] J. M. Girard, “CARMA: Software for continuous affect rating and media annotation,” *Journal of Open Research Software*, vol. 2, no. 1, p. e5, 2014. [Online]. Available: <https://carma.jmgirard.com>
- [21] K. Kercher, “Assessing Subjective Well-Being in the Old-Old: The PANAS as a Measure of Orthogonal Dimensions of Positive and Negative Affect,” *Research on Aging*, vol. 14, no. 2, pp. 131–168, 1992.
- [22] A. Mackinnon, A. F. Jorm, H. Christensen, A. E. Korten, P. A. Jacomb, and B. Rodgers, “A short form of the Positive and Negative Affect Schedule: evaluation of factorial validity and invariance across demographic variables in a community sample,” *Personality and Individual Differences*, vol. 27, no. 3, pp. 405–416, 1999.
- [23] D. ten Hove, T. D. Jorgensen, and L. A. van der Ark, “Updated guidelines on selecting an intraclass correlation coefficient for interrater reliability, with applications to incomplete observational designs,” *Psychological Methods*, Sep. 2022.
- [24] J. M. Girard, “varde: An R package for variance decomposition,” 2023. [Online]. Available: <https://github.com/jmgirard/varde>
- [25] T. K. Koo and M. Y. Li, “A guideline of selecting and reporting intraclass correlation coefficients for reliability research,” *Journal of Chiropractic Medicine*, vol. 15, no. 2, pp. 155–163, 2016.
- [26] R. P. McDonald, *Test theory: A unified approach*. Erlbaum, 1999.
- [27] D. B. Flora, “Your coefficient alpha is probably wrong, but which coefficient omega is right? A tutorial on using R to obtain better reliability estimates,” *Advances in Methods and Practices in Psychological Science*, vol. 3, no. 4, pp. 484–501, 2020.
- [28] Y. Rosseel, “lavaan: An R package for structural equation modeling,” *Journal of Statistical Software*, vol. 48, no. 2, pp. 1–36, 2012.
- [29] T. D. Jorgensen, S. Pornprasertmanit, A. M. Schoemann, and Y. Rosseel, “semTools: Useful tools for structural equation modeling,” 2022. [Online]. Available: <https://CRAN.R-project.org/package=semTools>
- [30] D. McNeish, “Thanks coefficient alpha, we’ll take it from here,” *Psychological Methods*, vol. 23, no. 3, pp. 412–433, 2017.
- [31] T. L. Kelley, “The applicability of the Spearman-Brown formula for the measurement of reliability,” *Journal of Educational Psychology*, vol. 16, no. 5, pp. 300–303, 1925.
- [32] R. Rosenthal, “Conducting judgment studies: Some methodological issues,” in *The new handbook of methods in nonverbal behavior research*, J. A. Harrigan, R. Rosenthal, and K. R. Scherer, Eds. Oxford University Press, 2005, pp. 199–234.
- [33] J. M. Girard, W.-S. Chu, L. A. Jeni, J. F. Cohn, F. De La Torre, and M. A. Sayette, “Sayette Group Formation Task (GFT) Spontaneous Facial Expression Database,” in *Proceedings of the 12th IEEE International Conference on Automatic Face & Gesture Recognition*, 2017, pp. 581–588.
- [34] J. McAuley, *Personalized machine learning*. Cambridge University Press, 2022.
- [35] D. Watson and A. Tellegen, “Toward a consensual model of mood,” *Psychological Bulletin*, vol. 98, no. 2, pp. 219–235, 1985.
- [36] J. M. Girard and A. G. C. Wright, “DARMA: Software for Dual Axis Rating and Media Annotation,” *Behavior Research Methods*, vol. 50, no. 3, pp. 902–909, 2018.
- [37] K. Fayn, S. Willemsen, R. Muralikrishnan, B. Castaño Manias, W. Menninghaus, and W. Schlotz, “Full throttle: Demonstrating the speed, accuracy, and validity of a new method for continuous two-dimensional self-report and annotation,” *Behavior Research Methods*, vol. 54, no. 1, pp. 350–364, 2022.
- [38] D. C. Rubin and J. M. Talarico, “A comparison of dimensional models of emotion: Evidence from emotions, prototypical events, autobiographical memories, and words,” *Memory*, vol. 17, no. 8, pp. 802–808, 2009.