# Multimodal Local-Global Ranking Fusion for Emotion Recognition

Paul Pu Liang Carnegie Mellon University Pittsburgh, PA, USA pliang@cs.cmu.edu Amir Zadeh Carnegie Mellon University Pittsburgh, PA, USA abagherz@cs.cmu.edu Louis-Philippe Morency Carnegie Mellon University Pittsburgh, PA, USA morency@cs.cmu.edu

# **ABSTRACT**

Emotion recognition is a core research area at the intersection of artificial intelligence and human communication analysis. It is a significant technical challenge since humans display their emotions through complex idiosyncratic combinations of the language, visual and acoustic modalities. In contrast to traditional multimodal fusion techniques, we approach emotion recognition from both direct person-independent and relative person-dependent perspectives. The direct person-independent perspective follows the conventional emotion recognition approach which directly infers absolute emotion labels from observed multimodal features. The relative person-dependent perspective approaches emotion recognition in a relative manner by comparing partial video segments to determine if there was an increase or decrease in emotional intensity. Our proposed model integrates these direct and relative prediction perspectives by dividing the emotion recognition task into three easier subtasks. The first subtask involves a multimodal local ranking of relative emotion intensities between two short segments of a video. The second subtask uses local rankings to infer global relative emotion ranks with a Bayesian ranking algorithm. The third subtask incorporates both direct predictions from observed multimodal behaviors and relative emotion ranks from local-global rankings for final emotion prediction. Our approach displays excellent performance on an audio-visual emotion recognition benchmark and improves over other algorithms for multimodal fusion.

# **CCS CONCEPTS**

• Computing methodologies  $\rightarrow$  Artificial Intelligence; Machine Learning;

# **KEYWORDS**

Multimodal Fusion; Neural Networks, Emotion Recognition

#### **ACM Reference Format:**

Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. 2018. Multimodal Local-Global Ranking Fusion for Emotion Recognition. In 2018 International Conference on Multimodal Interaction (ICMI '18), October 16–20, 2018, Boulder, CO, USA. ACM, New York, NY, USA, 5 pages. https://doi.org/10.1145/3242969. 3243019

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICMI '18, October 16–20, 2018, Boulder, CO, USA © 2018 Association for Computing Machinery. ACM ISBN 978-1-4503-5692-3/18/10. https://doi.org/10.1145/3242969.3243019

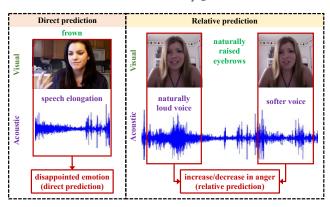


Figure 1: Our proposed Multimodal Local-Global Ranking Fusion (MLRF) model integrates both direct and relative prediction perspectives. The direct prediction perspective follows the conventional recognition approach which directly infers absolute emotion labels from observed multimodal features. The relative prediction perspective approaches emotion recognition in a relative manner by comparing partial video segments to determine changes in emotional intensity.

# 1 INTRODUCTION

Emotion recognition is a core research area at the intersection of artificial intelligence and human communication analysis. It has immense applications towards robotics [2, 15], dialog systems [20, 21], intelligent tutoring systems [16, 18, 28], and healthcare diagnosis [8]. Emotion recognition is multimodal in nature as humans utilize multiple communicative modalities in a structured fashion to convey emotions [4]. Two of these important modalities are acoustic and visual. In the acoustic modality, humans use prosody and various vocal expressions. In the visual modality, humans utilize facial expressions, hand gestures, and body language. Each modality is crucial when analyzing human emotions, making emotion recognition a challenging domain of artificial intelligence.

Some emotional expressions are almost universal and can be recognized directly from a video segment. For example, an open mouth with raised eyebrows and a loud voice is likely to be associated with surprise. These can be seen as *person-independent* behaviors and absolute emotions can be directly inferred from these behaviors (left panel of Figure 1). However, emotions are also expressed in a *person-dependent* fashion with idiosyncratic behaviors. In these cases, it may not be possible to directly estimate absolute emotion intensities. Instead, it would be easier to compare two video segments of the same person and judge whether there was a relative

change in emotion intensities (right panel of Figure 1). For example, a person could have naturally furrowed eyebrows and we should not always interpret this as a display of anger, but rather compare all nonverbal behaviors between two video segments to determine the relative changes in his displayed anger. From a psychological approach, research has also highlighted the advantages of using ordinal (relative) representations of human signals [6, 17, 26].

In this paper, we introduce the Multimodal Local-Global Ranking Fusion (MLRF) model which performs emotion recognition by integrating both *direct* prediction and *relative* prediction approaches. This is performed by dividing the emotion recognition task into three easier multimodal subtasks (Figure 2). The first subtask is the multimodal local ranking task. Given two short segments randomly selected from an entire video, the model is tasked with determining if there was an increase or decrease in the displayed emotion intensity. This task is often simpler than the direct emotion recognition problem since the model only needs to compare relative emotion ranks rather than compute the absolute intensities. The second subtask is the global ranking task, which uses the previous results of local rankings to infer relative global emotion ranks using a Bayesian skill rating algorithm [3, 10]. The third subtask involves direct-relative fusion of direct emotion predictions from observed multimodal behaviors with relative emotion ranks estimated from local-global ranking. This integration of direct and relative emotion predictions allows MLRF to model both person-independent and person-dependent behaviors for complete emotion recognition. We show that MLRF is suitable for multimodal tasks by performing experiments on an audio-visual emotion recognition benchmark. The proposed MLRF approach displays excellent performance over the baselines, improving over other algorithms for multimodal fusion.

# 2 RELATED WORK

Previous approaches in multimodal emotion recognition can be categorized as follows:

Non-temporal Models: These approaches simplify the temporal aspect by averaging information through time [1]. Fusion is performed by concatenating the multimodal inputs. However, these methods tend to overfit without discovering generalizable speaker-independent and speaker-dependent features [30]. More complex fusion methods learn separate models for each modality and combine the outputs [25]. However, simple decision voting is unable to discover the complex multimodal combinations involved in speaker-dependent features [13].

Temporal Models: Long Short-term Memory Networks (LSTMs) [9, 11, 24] have been extended for multimodal settings [22] and with binary gating mechanisms to remove noisy modalities [5]. More advanced models use memory mechanisms [32], low-rank approximations to tensor products [14], multiple attention stages [13] or assignments [33] or generative-discriminative objectives to learn factorized [27] or joint multimodal representations [19]. To our knowledge, our approach is the first to approach multimodal fusion with a neural local-global ranking fusion approach. The strength of our approach lies in approaching both speaker-independent and speaker-dependent features via direct and relative emotion predictions respectively. Algorithmically, our divide-and-conquer insight

simplifies the emotion recognition task into three easier multimodal subtasks and allows us to incorporate probabilistic structure as compared to entirely neural approaches.

Our work is also related to ranking algorithms. Bayesian ranking algorithms have been used in ranking the skills of players in Chess [7, 29] and online games [10]. Recently, ranking methods have been applied for facial expression intensity estimation [3]. To the best of our knowledge, we are the first to integrate relative measures from local-global ranking with direct predictions for emotion recognition. We also apply our approach to a multimodal setting where temporal information is primordial and there exist complex interactions between the acoustic and visual modalities.

# 3 MULTIMODAL LOCAL-GLOBAL RANKING FUSION

The Multimodal Local-Global Ranking Fusion (MLRF) model (Figure 2) aims to integrate both direct and relative emotion prediction approaches to model person-independent and person-dependent behaviors. This is achieved by subdividing the emotion recognition task into three easier subtasks: (1) *multimodal local ranking*, (2) *global ranking*, and (3) *direct-relative fusion*. Relative emotion rankings are inferred from the first two sub-tasks. The third subtask performs direct emotion predictions while at the same time integrating relative emotion rankings for final emotion recognition.

# 3.1 Problem Statement

Given a set of modalities  $\mathcal{M}$  each in the form of a temporal sequence with T time steps, we denote the data from modality  $m \in \mathcal{M}$  as  $\mathbf{x}^m = \langle \mathbf{x}_1^m, \mathbf{x}_2^m, \cdots, \mathbf{x}_T^m \rangle$ , where  $\mathbf{x}_t^m \in \mathbb{R}^{d_m}$  denotes the input of modality m at time t with dimensionality  $d_m$ . For a window size w, define  $\mathbf{x}_{t_w}^m = \langle \mathbf{x}_{t-w}^m, \cdots \mathbf{x}_t^m, \cdots \mathbf{x}_{t+w}^m \rangle$  as the short video segment centered around time t with a time window w. The goal is to estimate the sequence of emotion labels  $\mathbf{y} = \langle y_1, y_2, \cdots, y_T \rangle$ , where  $y_t \in \mathbb{R}$  is the emotion label at time t. In our experiments, the training dataset  $\mathcal{D}$  consists of n input-label pairs  $\mathcal{D} = \{\mathbf{x}_{(i)}^m, \mathbf{y}_{(i)} : m \in \mathcal{M}\}_{i=1}^n$ . The test dataset has a similar structure but with no overlapping speakers. The subscript (i) indicates the variable associated with the i-th video. We drop index (i) when it is clear from the context.

# 3.2 Multimodal Local Ranking

In the *multimodal local ranking task*, the model is presented with two short segments randomly selected within a video and is tasked with determining whether there was an increase or decrease in emotion intensity. This local ranking process is repeated multiple times for different segment pairs. The number of pairs is a hyperparameter and (J,K) denotes the set of all pairs. Given a pair  $(j,k) \in (J,K)$ ,  $\mathbf{x}_{j_w} = \{\mathbf{x}_{j_w}^m : m \in \mathcal{M}\}$  denotes a short segment integrating all multimodal features and likewise for  $\mathbf{x}_{k_w}$ . w is a hyper-parameter which defines the context over which local ranking is performed. The local rank  $r_{j,k} = \mathbb{I}[y_j > y_k]$  indicates whether there was an increase or decrease of emotion intensity between segments j and k. The multimodal local ranking dataset  $\mathcal{D}_{local}$  is therefore:

$$\mathcal{D}_{local} = \{ \{ \mathbf{x}_{j_{w(i)}}, \mathbf{x}_{k_{w(i)}}, r_{j,k(i)} \}_{j \in J_{(i)}, k \in K_{(i)}} \}_{i=1}^{n}$$
 (1)

This defines a binary classification problem over two short segments. This task is simpler than the direct emotion recognition

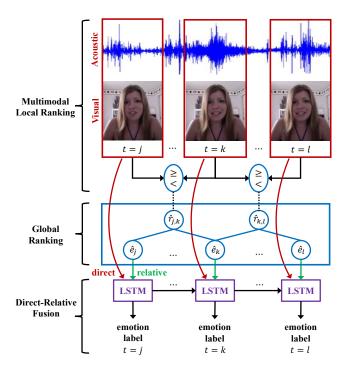


Figure 2: In MLRF, the emotion recognition task is divided into three easier subtasks: (1) multimodal local ranking involves a local ranking of emotion intensities between two short segments of a video, (2) global ranking uses the results of local rankings to infer global relative emotion ranks using a Bayesian ranking algorithm, and (3) direct-relative fusion integrates direct emotion predictions estimated from observed multimodal behaviors with relative emotion ranks from local-global rankings for final emotion recognition.

problem since the local ranking model only needs to compare the differences in emotion intensities rather than compute the exact emotion intensities themselves.

To solve the multimodal local ranking problem, we define an estimator  $f_{local}$  with parameters  $\Theta_{local}$ :

$$\hat{r}_{j,k} = f_{local}(\mathbf{x}_{j_w}, \mathbf{x}_{k_w}; \Theta_{local})$$
 (2)

Solving this problem involves multimodal fusion since  $\mathbf{x}_{j_w}$  and  $\mathbf{x}_{k_w}$  represent data from the input modalities  $\mathcal{M}$ . In order to solve for parameters  $\Theta_{local}$ , we minimize the empirical measure of the categorical cross-entropy between the target local ranks  $r_{j,k}$  and our estimated local ranks  $\hat{r}_{i,k}$ :

$$\hat{\mathcal{L}}_{local} = \frac{1}{|\mathcal{D}_{local}|} \sum_{(\mathbf{x}_{j_{w}}, \mathbf{x}_{k_{w}}, r_{j,k}) \in \mathcal{D}_{local}} r_{j,k} \log \hat{r}_{j,k}$$
(3)

In practice, we parametrize  $f_{local}$  using an LSTM [11] which takes as input the differences of multimodal tensors:

$$\mathbf{x}_{local} = \bigoplus_{m \in \mathcal{M}} \mathbf{x}_{k_w}^m - \bigoplus_{m \in \mathcal{M}} \mathbf{x}_{j_w}^m \tag{4}$$

where  $\oplus$  denotes tensor concatenation of all unimodal feature vectors  $\mathbf{x}^m$ .  $\mathbf{x}_{local}$  is the input sequence to the LSTM that performs multimodal local ranking. A neural network classification layer on the final LSTM output  $\mathbf{h}_T$  is used to estimate local ranks  $\hat{r}_{i,k}$ .

In contrast to [3] who developed a model for images only, our problem involves multimodal video segments. In our models, local comparisons of emotion intensities are performed over a time window *w*. Emotion intensities often require more than still frames, especially when including the acoustic modality. Furthermore, human communicative behaviors can be asynchronous and a longer time window is required to track changes in both visual and acoustic behaviors displayed by the person.

# 3.3 Global Ranking

The second task is the *global ranking* task, which uses the previous results of local rankings to infer global emotion ranks using a Bayesian skill rating algorithm [10]. The algorithm will infer global emotion ranks  $e_t$  at each time step t of the multimodal video. These ranks are initially sampled from a prior distribution  $\mathcal{N}(\mu_t, \sigma_t^2)$ . The algorithm models global emotion ranks  $e_t$  as hidden variables that are not directly observed from the data. What we do observe are local rankings  $r_{j,k}$  between two time segments j and k. By our definition of local ranks, the conditional probabilities of local ranks  $p(r_{j,k}=1|e_j,e_k)$  will be equal to  $p(e_j>e_k)$ , the probability of rank  $e_j$  being larger than  $e_k$ . Given estimated ranks  $\hat{r}_{j,k}$ , we estimate global relative emotion ranks  $\hat{e}_t$  using a ranking algorithm which involves message passing over factor graph models. A detailed treatment is presented in [3, 10]. After all iterations of global ranking, we obtain estimated global emotion ranks  $\hat{e}_t$  for all video segments.

## 3.4 Direct-Relative Fusion

The third task involves direct-relative fusion. The global emotion ranks  $\hat{\mathbf{e}} = \langle \hat{e}_1, \cdots, \hat{e}_T \rangle$  are incorporated with the raw multimodal inputs  $\mathbf{x}^m$  to estimate final emotion intensities. This allows us to perform direct estimation of the absolute emotion intensities from multimodal data while at the same time integrating relative emotion ranks from local-global rankings.

The integration of direct and relative predictions is performed by learning a function  $f_{fusion}$  with parameters  $\Theta_{fusion}$ :

$$\hat{\mathbf{y}} = f_{fusion}(\mathbf{x}, \hat{\mathbf{e}}; \Theta_{fusion})$$
 (5)

where  $\hat{\mathbf{y}}$  are the predicted emotion labels. We solve for  $\Theta_{fusion}$  by minimizing an empirical measure of the loss between  $\hat{\mathbf{y}}$  and  $\mathbf{y}$ :

$$\hat{\mathcal{L}}_{fusion} = \frac{1}{n} \sum_{i=1}^{n} \ell(\hat{\mathbf{y}}_{(i)}, \mathbf{y}_{(i)})$$
 (6)

where  $\ell(\hat{y}, y)$  is a loss function over two segments of emotion intensities. In practice, we parametrize  $f_{fusion}$  with an LSTM on the concatenated multimodal data and global emotion ranks:

$$\mathbf{x}_{fusion} = \left(\bigoplus_{m \in \mathcal{M}} \mathbf{x}^m\right) \bigoplus \hat{\mathbf{e}}$$
 (7)

where  $\oplus$  denotes tensor concatenation. Using an LSTM on  $\mathbf{x}_{fusion}$  allows us to capture the temporal dependencies across multimodal time series data and label emotion intensities. A time-distributed neural network regression layer is used on the LSTM outputs  $\mathbf{h}_t$  to estimate final emotion intensities  $\hat{y}_t$ . The global emotion ranks  $\hat{\mathbf{e}}$  can be incorporated with multimodal data  $\mathbf{x}^m$  using any multimodal fusion method. As a result, our approach represents a generalizable framework to integrate relative emotion intensities into a variety of multimodal fusion models.

Dataset	AVEC16	
Task	Arousal	Valence
Metric	CCC	CCC
EF-(-/S/B/SB)LSTM [9, 11, 24]	0.4327	0.4667
Gated-LSTM [31]	0.3210	0.4667
MV-LSTM, view-specific [22]	0.4530	0.4431
MV-LSTM, coupled [22]	0.4300	0.4477
MV-LSTM, hybrid [22]	0.4729	0.4924
MV-LSTM, fully connected [22]	0.4293	0.4896
MLRF-500	0.4732	0.5063
MLRF-1000	0.5049	0.5432
Improvement over baselines	↑ 0.032	↑ 0.0508

Table 1: Multimodal emotion recognition results on AVEC16 dataset. The best results are highlighted in bold and green indicates the improvement over baselines. The MLRF outperforms the baselines across all evaluation metrics.

## 4 EXPERIMENTS

#### 4.1 Dataset

We use the AVEC16 dataset (RECOLA) [23] for audio-visual emotion recognition. AVEC16 consists of 9 training videos, 8 are used to optimize parameters (our training set) and 1 is held out to tune hyperparameters (our validation set). The 9 validation videos are used as our test set to compare each method (test videos are not publicly released). We use the provided appearance and geometric visual features and acoustic features. Each video has 7501 time steps after alignment between the modalities and is labeled continuously for arousal and valence at every time step. The metric used is the concordance correlation coefficient (CCC) [12].

# 4.2 Baseline Models

We compare to the following: **EF-LSTM** (Early Fusion LSTM) uses a single LSTM [11] on concatenated multimodal inputs. We also implement the **EF-SLSTM** (stacked) [9], **EF-BLSTM** (bidirectional) [24] and **EF-SBLSTM** (stacked bidirectional) versions and report the best result. **GF-LSTM** (Gated Fusion LSTM) [31] extends the EF-LSTM by assigning an LSTM to each modality and combining the final LSTM outputs with a gated attention fusion layer for final prediction. **MV-LSTM** (Multi-View LSTM) [22] allocates regions inside a LSTM to different modalities using parameters  $\alpha$  and  $\beta$ . We experiment with the view-specific  $(0 < \alpha \le 1, \beta = 0)$ , coupled  $(\alpha = 0, 0 < \beta \le 1)$ , hybrid  $(0 < \alpha < 1, 0 < \beta < 1)$  and fully connected  $(\alpha = 1, \beta = 1)$  topologies as well. Our model is indicated by **MLRF-k** where k is the number of local comparison pairs. We set the default window size for local ranking as w = 200.

## 4.3 Results on Emotion Recognition

MLRF achieves better results on arousal and valence regression as compared to the baselines (Table 1). Our results show that incorporating local-global ranking estimates into simple models (EF-LSTM in our experiments) prove more effective than engineering complex neural architectures such as the Gated-LSTM [31] and the MV-LSTM [22]. Furthermore, although the videos contain more than 7500 time steps, simply sampling 500-1000 local ranking pairs per video significantly improved final performance. As a result, incorporating relative emotion intensities via MLRF is an effective method without adding significant computational complexity.

Dataset	AVEC16	
Task	Arousal	Valence
Metric	CCC	CCC
MLRF-500 $w = 10$	0.4165	0.2377
MLRF-500 $w = 50$	0.4168	0.4175
MLRF-500 $w = 100$	0.4196	0.4340
MLRF-500 $w = 200^{\circ}$	0.4732	0.5063

Table 2: Increasing the window size w improves performance. We observed a similar trend for MLRF-1000.

Dataset	AVEC16	
Task	Arousal	Valence
Metric	CCC	CCC
MLRF-500 direct predictions only	0.4327	0.4667
MLRF-500 relative predictions only	0.3646	0.0402
MLRF-500	0.4732	0.5063
MLRF-1000 direct predictions only	0.4327	0.4667
MLRF-1000 relative predictions only	0.4297	0.0846
MLRF-1000	0.5049	0.5432

Table 3: Ablation studies: incorporating both direct and relative emotion predictions is crucial for performance.

# 4.4 Discussion

**Effect of Number of Local Comparison Pairs:** Table 1 shows that performance increases as the number of sampled local comparison pairs increases. More observations of local ranks  $\hat{r}_{j,k}$  improves our estimates of global emotion ranks  $\hat{e}_t$ , which in turn provide better relative emotion intensities for emotion recognition.

**Effect of Window Size:** From Table 2, we observe that increasing the window size w for multimodal local ranking is important. This supports the fact that human communicative behaviors are asynchronous and a longer time window is required to track changes in visual and acoustic behaviors displayed by the speaker.

**Effect of Direct and Relative Approaches:** We found that fusing direct emotion predictions from observed multimodal behaviors with relative emotion predictions from local-global ranking estimates was crucial (Table 3). Therefore, integrating both direct person-independent and relative person-dependent approaches is important for emotion recognition.

## 5 CONCLUSION

This paper approached multimodal emotion recognition from both direct person-independent and relative person-dependent perspectives. Our proposed Multimodal Local-Global Ranking Fusion (MLRF) model integrates direct and relative predictions by dividing emotion recognition into three easier subtasks: *multimodal local ranking*, *global ranking* and *direct-relative fusion*. Our experiments showed that MLRF displays excellent performance on multimodal tasks. Therefore, incorporating direct emotion predictions from multimodal behaviors and relative emotion ranks from local-global rankings is a promising direction for multimodal machine learning.

# **6 ACKNOWLEDGEMENTS**

This material is based upon work partially supported by Samsung. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of Samsung, and no official endorsement should be inferred.

#### REFERENCES

- Harika Abburi, Rajendra Prasath, Manish Shrivastava, and Suryakanth V Gangashetty. 2016. Multimodal Sentiment Analysis Using Deep Neural Networks. In International Conference on Mining Intelligence and Knowledge Exploration. Springer, 58–65.
- [2] Fernando Alonso-Martin, Maria Malfaz, Joao Sequeira, Javier F. Gorostiza, and Miguel A. Salichs. 2013. A Multimodal Emotion Detection System during Human-Robot Interaction. Sensors 13, 11 (2013), 15549–15581. https://doi.org/10.3390/ s131115549
- [3] T. Baltrušaitis, L. Li, and L. P. Morency. 2017. Local-global ranking for facial expression intensity estimation. In 2017 Seventh International Conference on Affective Computing and Intelligent Interaction (ACII). 111–118. https://doi.org/10. 1109/ACII.2017.8273587
- [4] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette Chang, Sungbok Lee, and Shrikanth S. Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Journal of Language Resources and Evaluation* 42, 4 (dec 2008), 335–359. https://doi.org/10.1007/s10579-008-9076-6
- [5] Minghai Chen, Sen Wang, Paul Pu Liang, Tadas Baltrušaitis, Amir Zadeh, and Louis-Philippe Morency. 2017. Multimodal Sentiment Analysis with Word-level Fusion and Reinforcement Learning. In Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI 2017). ACM, New York, NY, USA, 163–171. https://doi.org/10.1145/3136755.3136801
- [6] R. Elliott, Z. Ágnew, and J. F. W. Deakin. Medial orbitofrontal cortex codes relative rather than absolute value of financial rewards in humans. European Journal of Neuroscience 27, 9 (????), 2213–2218. https://doi.org/10.1111/j.1460-9568. 2008.06202.x arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1460-9568.2008.06202.x
- [7] Arpad E. Elo. 1978. The rating of chessplayers, past and present. Arco Pub., New York. http://www.amazon.com/Rating-Chess-Players-Past-Present/dp/ 0668047216
- [8] C. A. Frantzidis, C. Bratsas, M. A. Klados, E. Konstantinidis, C. D. Lithari, A. B. Vivas, C. L. Papadelis, E. Kaldoudi, C. Pappas, and P. D. Bamidis. 2010. On the Classification of Emotional Biosignals Evoked While Viewing Affective Pictures: An Integrated Data-Mining-Based Approach for Healthcare Applications. IEEE Transactions on Information Technology in Biomedicine 14, 2 (March 2010), 309–318. https://doi.org/10.1109/TITB.2009.2038481
- [9] A. Graves, A. r. Mohamed, and G. Hinton. 2013. Speech recognition with deep recurrent neural networks. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing. 6645–6649. https://doi.org/10.1109/ICASSP.2013. 6638047
- [10] Ralf Herbrich, Tom Minka, and Thore Graepel. 2007. TrueSkill™: A Bayesian Skill Rating System. In Advances in Neural Information Processing Systems 19, B. Schölkopf, J. C. Platt, and T. Hoffman (Eds.). MIT Press, 569–576. http://papers. nips.cc/paper/3079-trueskilltm-a-bayesian-skill-rating-system.pdf
- [11] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. Neural computation 9, 8 (1997), 1735–1780.
- [12] Lawrence I-Kuei Lin. 1989. A Concordance Correlation-Coefficient To Evaluate Reproducibility. 45 (04 1989), 255–68.
- [13] Paul Pu Liang, Ziyin Liu, Amir Zadeh, and Louis-Philippe Morency. 2018. Multi-modal Language Analysis with Recurrent Multistage Fusion. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2018).
- [14] Zhun Liu, Ying Shen, Varun Bharadhwaj Lakshminarasimhan, Paul Pu Liang, AmirAli Bagher Zadeh, and Louis-Philippe Morency. 2018. Efficient Low-rank Multimodal Fusion With Modality-Specific Factors. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, 2247–2256. http://aclweb. org/anthology/P18-1209
- [15] Z. Liu, M. Wu, W. Cao, L. Chen, J. Xu, R. Zhang, M. Zhou, and J. Mao. 2017. A facial expression emotion recognition based human-robot interaction system. *IEEE/CAA Journal of Automatica Sinica* 4, 4 (2017), 668–676. https://doi.org/10. 1109/JAS.2017.7510622
- [16] Mehdi Malekzadeh, Mumtaz Begum Mustafa, and Adel Lahsasna. 2015. A Review of Emotion Regulation in Intelligent Tutoring Systems. *Journal of Educational Technology and Society* 18, 4 (2015), 435–445. http://www.jstor.org/stable/jeductechsoci.18.4.435
- [17] George Miller. 1956. The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information. (1956), 81–97 pages. http://cogprints.org/730/ One of the 100 most influential papers in cognitive science: http://cogsci.umn.edu/millennium/final.html.
- [18] Sintija Petrovica, Alla Anohina-Naumeca, and HazÄśm Kemal Ekenel. 2017. Emotion Recognition in Affective Tutoring Systems: Collection of Ground-truth Data. Procedia Computer Science 104 (2017), 437 – 444. https://doi.org/10.1016/j. procs.2017.01.157 ICTE 2016, Riga Technical University, Latvia.
- [19] Hai Pham, Thomas Manzini, Paul Pu Liang, and Barnabas Poczos. 2018. Seq2Seq2Sentiment: Multimodal Sequence to Sequence Models for Sentiment

- Analysis. In Proceedings of Grand Challenge and Workshop on Human Multi-modal Language (Challenge-HML). Association for Computational Linguistics, Melbourne, Australia, 53–63. http://www.aclweb.org/anthology/W18-3308
- [20] Johannes Pittermann, Angela Pittermann, and Wolfgang Minker. 2009. Handling Emotions in Human-Computer Dialogues (1st ed.). Springer Publishing Company, Incorporated.
- [21] Johannes Pittermann, Angela Pittermann, and Wolfgang Minker. 2010. Emotion recognition and adaptation in spoken dialogue systems. *International Jour*nal of Speech Technology 13, 1 (01 Mar 2010), 49–60. https://doi.org/10.1007/ s10772-010-9068-v
- [22] Shyam Sundar Rajagopalan, Louis-Philippe Morency, Tadas Baltrušaitis, and Roland Goecke. 2016. Extending long short-term memory for multi-view structured learning. In European Conference on Computer Vision.
- [23] Fabien Ringeval, Andreas Sonderegger, Jürgen S. Sauer, and Denis Lalanne. 2013. Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions.. In FG. IEEE Computer Society, 1–8.
- [24] M. Schuster and K.K. Paliwal. 1997. Bidirectional Recurrent Neural Networks. Trans. Sig. Proc. 45, 11 (Nov. 1997), 2673–2681. https://doi.org/10.1109/78.650093
- [25] Cees G. M. Snoek, Marcel Worring, and Arnold W. M. Smeulders. 2005. Early Versus Late Fusion in Semantic Video Analysis. In Proceedings of the 13th Annual ACM International Conference on Multimedia (MULTIMEDIA '05). ACM, New York, NY, USA, 399–402. https://doi.org/10.1145/1101149.1101236
- [26] Neil Stewart, Gordon D. A. Brown, and Nick Chater. 2005. Absolute identification by relative judgment. Psychological review 112 4 (2005), 881–911.
- [27] Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2018. Learning Factorized Multimodal Representations. arXiv preprint arXiv:1806.06176 (2018).
- [28] Alexandria K. Vail, Joseph F. Grafsgaard, Kristy Elizabeth Boyer, Eric N. Wiebe, and James C. Lester. 2016. Gender Differences in Facial Expressions of Affect During Learning. In Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization (UMAP '16). ACM, New York, NY, USA, 65–73. https://doi.org/10.1145/2930238.2930257
- [29] N. Veček, M. Črepinšek, M. Mernik, and D. Hrnčič. 2014. A comparison between different chess rating systems for ranking evolutionary algorithms. In 2014 Federated Conference on Computer Science and Information Systems. 511–518. https://doi.org/10.15439/2014F33
- [30] Chang Xu, Dacheng Tao, and Chao Xu. 2013. A survey on multi-view learningarXiv preprint arXiv:1304.5634 (2013).
- [31] Hongliang Yu, Liangke Gui, Michael Madaio, Amy Ogan, Justine Cassell, and Louis-Philippe Morency. 2017. Temporally Selective Attention Model for Social and Affective State Recognition in Multimedia Content. In Proceedings of the 2017 ACM on Multimedia Conference (MM '17). ACM, New York, NY, USA, 1743–1751. https://doi.org/10.1145/3123266.3123413
- [32] Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018. Memory Fusion Network for Multi-view Sequential Learning. Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (2018).
- [33] Amir Zadeh, Paul Pu Liang, Soujanya Poria, Prateek Vij, Erik Cambria, and Louis-Philippe Morency. 2018. Multi-attention recurrent network for human communication comprehension. Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (2018).