

MMLATCH: BOTTOM-UP TOP-DOWN FUSION FOR MULTIMODAL SENTIMENT ANALYSIS

Georgios Paraskevopoulos^{†,*} Efthymios Georgiou^{†,*} Alexandros Potamianos^{†,‡}

[†] National Technical University of Athens, Athens, Greece

* Institute for Language and Speech Processing, Athena Research Center, Athens, Greece

[‡] Behavioral Signal Technologies, Los Angeles, CA, USA

ABSTRACT

Current deep learning approaches for multimodal fusion rely on bottom-up fusion of high and mid-level latent modality representations (late/mid fusion) or low level sensory inputs (early fusion). Models of human perception highlight the importance of top-down fusion, where high-level representations affect the way sensory inputs are perceived, i.e. cognition affects perception. These top-down interactions are not captured in current deep learning models. In this work we propose a neural architecture that captures top-down cross-modal interactions, using a feedback mechanism in the forward pass during network training. The proposed mechanism extracts high-level representations for each modality and uses these representations to mask the sensory inputs, allowing the model to perform top-down feature masking. We apply the proposed model for multimodal sentiment recognition on CMU-MOSEI. Our method shows consistent improvements over the well established MulT and over our strong late fusion baseline, achieving state-of-the-art results.

Index Terms— multimodal, fusion, sentiment, feedback

1. INTRODUCTION

Multimodal processing aims to model interactions between inputs that come from different sources in real world tasks. Multimodality can open ways to develop novel applications (e.g. Image Captioning, Visual Question Answering [1, 2] etc.) or boost performance in traditionally unimodal applications (e.g. Machine Translation [3], Speech Recognition [4, 5] etc.). Moreover, modern advances in neuroscience and psychology hint that multi-sensory inputs are crucial for cognitive functions [6], even since infancy [7]. Thus, modeling and understanding multimodal interactions can open avenues to develop smarter agents, inspired by the human brain.

Feedback loops have been shown to exist in the human brain, e.g. in the case of vocal production [8] or visual-motor coordination [9]. Human perception has been traditionally modelled as a linear (bottom-up) process (e.g. reflected light is captured by the eye, processed in the prefrontal visual cortex, then the posterior visual cortex etc.). Recent studies have

highlighted that this model may be too simplistic and that high level cognition may affect low-level visual [10, 11] or audio [12] perception. For example, studies state that perception may be affected by an individual’s long-term memory [13], emotions [14] and physical state [15]. Researchers have also tried to identify brain circuits that allow for this interplay [16]. While scientists still debate on this subject [17], such works offer strong motivation to explore if artificial neural networks can benefit from multimodal top-down modeling.

Early works on multimodal machine learning use binary decision trees [18] and ensembles of Support Vector Machines [19]. Modeling contextual information is addressed in [20, 21, 22] using Recurrent Neural Networks (RNNs), while Poria et al. [23] use Convolutional Neural Networks (CNNs). For a detailed review we refer to Baltruvsaitis et al. [24]. Later works use Kronecker product between late representations [25, 26], while others investigate architectures with neural memory-like modules [27, 28]. Hierarchical attention mechanisms [29], as well as hierarchical fusion [30] have been also proposed. Pham et al. [31] learn cyclic cross-modal mappings, Sun et al. [32] propose Deep Canonical Correlation Analysis (DCCA) for jointly learning representations. Multitask learning has been also investigated [33] in the multimodal context. Transformers [34] have been applied to and extended for multimodal tasks [35, 36, 37, 38]. Wang et al. [39] shift word representations based on non-verbal information. [40] propose a fusion gating mechanism. [41] use capsule networks [42] to weight input modalities and create distinct representations for input samples.

In this work we propose MMLatch, a neural network module that uses representations from higher levels of the architecture to create top-down masks for the low level input features. The masks are created by a set of feedback connections. The module is integrated in a strong late fusion baseline based on LSTM [43] encoders and cross-modal attention. Our key contribution is the modeling of interactions between high-level representations extracted by the network and low-level input features, using an end to end framework. We integrate MMLatch with RNNs, but it can be adapted for other architectures (e.g. Transformers). Incorporating

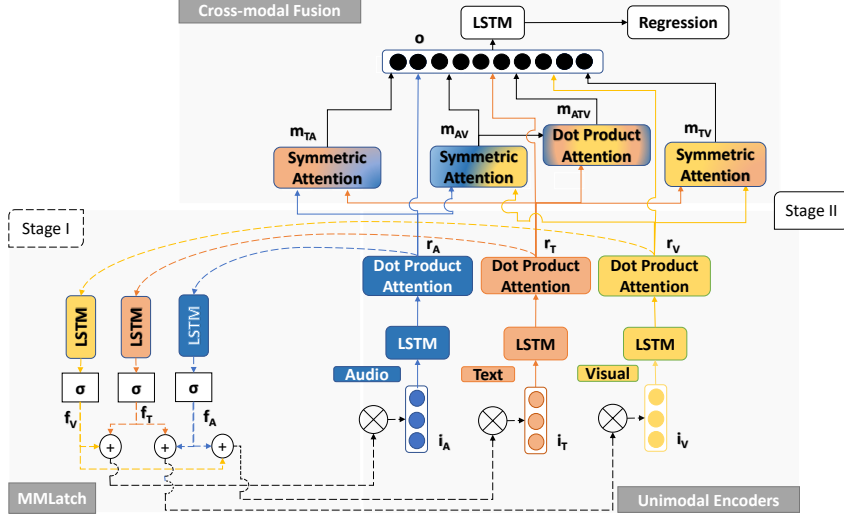


Fig. 1. Architecture overview of three high-level modules, composing the overall system: Unimodal encoders, Cross-modal fusion and MMLatch. Solid lines indicate the feedforward connections (bottom-up processing), while dashed lines indicate feedback connections (top-down processing). Colors indicate different modalities (Blue: Audio, Orange: Text, Yellow: Visual)

top-down modeling shows consistent improvements over our strong baseline, yielding state-of-the-art results for sentiment analysis on CMU-MOSEI. Qualitative analysis of learned top-down masks can add interpretability in multimodal architectures. Our code will be made available as open source.

2. PROPOSED METHOD

Fig. 1 illustrates an overview of the system architecture. The baseline system consists of a set of unimodal encoders and a cross-modal attention fusion network, that extracts fused feature vectors for regression on the sentiment values. We integrate top-down information by augmenting the baseline system with a set of feedback connections that create cross-modal, top-down feature masks.

Unimodal Encoders: Input features i_A, i_T, i_V for each modality are encoded using three LSTMs. The hidden states of each LSTM are then passed through a Dot Product self-attention mechanism to produce the unimodal representations r_A, r_T, r_V , where A, T, V are the audio, text and visual modalities respectively.

Cross-modal Fusion: The encoded unimodal representations are fed into a cross-modal fusion network, that uses a set of attention mechanisms to capture cross-modal interactions. The core component of this subsystem is the symmetric attention mechanism, inspired by Lu et al. [36]. If we consider modality indicators $k, l \in \{A, V, T\}$, $k \neq l$, $r_k, r_l \in \mathbb{R}^{B \times N \times d}$ the input modality representations, we can construct keys $K_l = W_l^K r_l$, queries $Q_k = W_k^Q r_k$ and values $V_l = W_l^V r_l$ using learnable projection matrices $W_{\{k,l\}}^{\{K,Q,V\}}$, and we can define a cross-modal attention layer as:

$$a_{kl} = s \left(\frac{K_l^T Q_k}{\sqrt{d}} \right) V_l + r_k, \quad (1)$$

where s is the softmax operation and B, N, d are the batch size, sequence length and hidden size respectively. For the symmetric attention we sum the two cross-modal attentions:

$$m_{kl} = a_{kl} + a_{lk}, \quad (2)$$

In the fusion subsystem we use three symmetric attention mechanisms to produce m_{TA}, m_{TV} and m_{AV} . Additionally we create a_{AVT} using a cross-modal attention mechanism (Eq. (1)) with inputs m_{AV} and r_T . These crossmodal representations are concatenated (\parallel), along with the unimodal representations m_A, m_T, m_V to produce the fused feature vector $o \in \mathbb{R}^{B \times N \times 7d}$ in Eq. (3).

$$o = r_A \parallel r_T \parallel r_V \parallel a_{AVT} \parallel m_{AV} \parallel m_{TV} \parallel m_{TA} \quad (3)$$

We then feed o into a LSTM and the last hidden state is used for regression. The baseline system consists of the unimodal encoders followed by the cross-modal fusion network. **Top-down fusion:** We integrate top-down information by augmenting the baseline system with MMLatch, i.e. a set of feedback connections composing of three LSTMs followed by sigmoid activations σ . The inputs to these LSTMs are r_A, r_T, r_V as they come out of the unimodal encoders. Feedback LSTMs produce hidden states h_A, h_T, h_V . The feedback masks f_A, f_T, f_V are produced by applying a sigmoid activation on the hidden states $f_k = \sigma(h_k)$, $k \in \{A, T, V\}$ and then applied to the input features i_A, i_T, i_V using element-wise multiplication \odot , as:

Model / Metric	Acc@7	Acc@2	F1@2	MAE	Corr
RAVEN [39] *	50.0	79.1	79.5	0.614	0.662
MCTN [31] *	49.6	79.8	80.6	0.609	0.670
Multimodal Routing [41]	51.6	81.7	81.8	-	-
MuT [35]	51.8	82.5	82.3	0.580	0.703
Baseline (ours)	51.3 ± 0.7	81.9 ± 0.7	82.2 ± 0.6	0.593 ± 0.005	0.695 ± 0.004
Baseline + MMLatch average (ours)	52.0 ± 0.2	82.4 ± 0.3	82.5 ± 0.3	0.585 ± 0.002	0.700 ± 0.004
Baseline + MMLatch best (ours)	52.1	82.8	82.9	0.582	0.704

Table 1. Results on CMU-MOSEI for MMLatch. Models indicated with * are reproduced for CMU-MOSEI by Tsai et al. [35]. In row “MMLatch average” we include results averaged over five runs. Since other works do not report standard deviation, we also include row “MMLatch best”, where we report the best of the five runs (lowest error).

$$\tilde{i}_k = \frac{1}{2}(f_j + f_l) \odot i_k \quad (4)$$

where $j, k, l \in \{A, V, T\}$, $k \neq l \neq m$.

Eq. (4) describes how the feedback masks for two modalities are applied to the input features of the third. For example, consider the case where we mask visual input features using the (halved) sum of text and audio feedback masks. If a visual feature is important for both audio and text the value of the resulting mask will be close to 1. If it is important for only one other modality the value will be close to 0.5, while if it is irrelevant for text and audio the value will be close to 0. Thus, a feature is enhanced or attenuated based on its overall importance for cross-modal representations.

This pipeline is implemented as a two-stage computation. During the first stage we use the unimodal encoders and MMLatch to produce the feedback masks f_A, f_T, f_V and apply them to the input features using Eq. (4). During the second stage we pass the masked features $\tilde{i}_A, \tilde{i}_T, \tilde{i}_V$ through the unimodal encoders and the cross-modal fusion module and use the fused representations for regression.

3. EXPERIMENTAL SETUP

We use CMU-MOSEI sentiment analysis dataset [28] for our experiments. The dataset contains 23,454 YouTube video clips of movie reviews accompanied by human annotations for sentiment scores from -3 (strongly negative) to 3 (strongly positive) and emotion annotations. Audio sequences are sampled at 20Hz and then 74 COVAREP features are extracted. Visual sequences are sampled at 15Hz and represented using Facet features. Video transcriptions are segmented in words and represented using GloVe. All sequences are word-aligned using P2FA. Standard train, validation and test splits are provided.

For all our experiments we use bidirectional LSTMs with hidden size 100. LSTMs are bidirectional and forward and backward passes are summed. All projection sizes for the attention modules are set to 100. We use dropout 0.2. We use Adam [44] with learning rate 0.0005 and halve the learning rate if the validation loss does not decrease for 2 epochs. We

use early stopping on the validation loss (patience 10 epochs). During Stage I of each training step we disable gradients for the unimodal encoders. Models are trained for regression on sentiment values using Mean Absolute Error (MAE) loss. We use standard evaluation metrics: 7-class, 5-class accuracy (i.e. classification in $\mathbb{Z} \cap [-3, 3]$, $\mathbb{Z} \cap [-2, 2]$), binary accuracy and F1-score (negative in $[-3, 0]$, positive in $(0, 3]$), MAE and correlation between model and human predictions. For fair comparison we compare with methods in the literature that use Glove text features, COVAREP audio features and FACET visual features.

4. EXPERIMENTS

Table 1 shows the results for sentiment analysis on CMU-MOSEI. The Baseline row refers to our late-fusion baseline described in Section 2, which achieves competitive to the state-of-the-art performance. Incorporating MMLatch into the baseline consistently improves performance and specifically, almost 1.0% over the binary accuracy and 0.8% over the seven class accuracy. Moreover, we observe lower deviation, w.r.t. the baseline, across experiments, indicating that top-down feedback can stabilize training. Compared to state-of-the-art we achieve better performance for 7-class accuracy and binary F1 metrics in our five run experiments. Since, prior works do not report average results over multiple runs so we also report results for the best (mean absolute error) out of five runs in the last row of Table 1, showing improvements across metrics over the best runs of the other methods.

In Table 2 we evaluate MMLatch with different multimodal encoders and different feedback types. The first three rows show the effect of using different feedback types. Specifically, first row shows our baseline performance (no feedback). For the second row we add feedback connections, but instead of using LSTMs in the feedback loop (Stage I in Fig. 1), we use a simple feed-forward layer. The last row shows performance when we include LSTMs in the feedback loop. We observe that, while the inclusion of top-down feedback, using a simple projection layer results to a small performance boost, when we include an LSTM in the feedback loop we get significant improvements. This shows that

Multimodal Encoder	Feedback Type	Acc@7	Acc@2	F1@2	MAE	Corr
Baseline	-	51.3 ± 0.7	81.9 ± 0.7	82.2 ± 0.6	0.593 ± 0.005	0.695 ± 0.004
Baseline	MMLatch (no LSTM)	51.48 ± 0.41	82.07 ± 0.47	82.29 ± 0.39	0.592 ± 0.002	0.692 ± 0.003
Baseline	MMLatch	52.0 ± 0.2	82.4 ± 0.3	82.5 ± 0.3	0.585 ± 0.002	0.700 ± 0.004
MuT [†]	-	47.91 ± 1.13	80.35 ± 0.36	80.54 ± 0.52	0.643 ± 0.01	0.648 ± 0.02
MuT [†]	MMLatch	49.04 ± 0.45	80.65 ± 0.43	81.07 ± 0.38	0.627 ± 0.004	0.665 ± 0.003

Table 2. Results on CMU-MOSEI when combining top-down feedback with different multimodal encoder networks. MuT with [†] is reproduced by us. We report results, averaged over five runs, along with standard deviations.

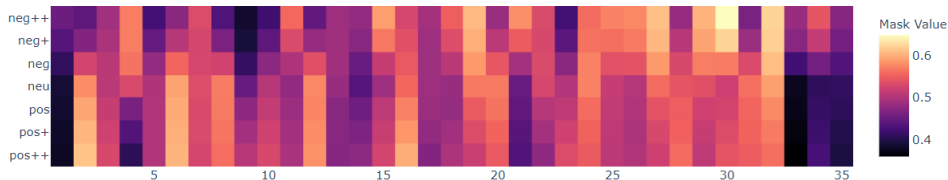


Fig. 2. Averaged top-down mask values for Facet features over all test samples across seven sentiment classes. neg++ indicates a sentiment score ≈ -3 , neg+ ≈ -2 , neg ≈ -1 , neu ≈ 0 , pos ≈ 1 , pos+ ≈ 2 and pos++ ≈ 3 .

choosing an appropriate mapping from high-level representations to low-level features in the feedback loop is important.

For the last two rows of Table 2 we integrate MMLatch with MuT architecture¹ [35]. Specifically, we use MMLatch, as shown in Fig. 1 and swap the baseline architecture (unimodal encoders and cross-modal fusion) with MuT. We use a 4-layer Transformer model with the same hyperparameter set and feature set described in the original paper [35]. The output of the fourth (final) layer is used by MMLatch to mask the input features. First, we notice a performance gap between our reproduced results and the ones reported in the original paper (fourth row of Table 2). Other works [45, 46] have reported similar observations. We observe that the integration of MMLatch with MuT yields significant performance improvements across metrics. Furthermore, similarly to Table 1, we observe that the inclusion of MMLatch reduces standard deviation across metrics. Overall, we observe that the inclusion of MMLatch results to performance improvements for both our baseline model and MuT with no additional tuning, indicating that top-down feedback can provide stronger multimodal representations.

Fig. 2 shows a heatmap of the average mask values $\frac{1}{2}(f_T + f_A)$. This mask is applied to the input visual features i_V , i.e. 35 Facet features. The average mask values range from 0.36 to 0.65 and depicted across 7 sentiment classes. Some features are attenuated or enhanced across all classes (e.g. features 1 or 32). Interestingly, some features are attenuated for some classes and enhanced for others (e.g. feature 2). More importantly this transition is smooth, i.e. mask values change almost monotonically as the sentiment value increases from -3 to $+3$, indicating well-behaved training of MMLatch. We observe the same for Covarep masks.

¹We use the original code in this GitHub Link

5. CONCLUSIONS

We introduce MMLatch, a feedback module that allows modeling top-down cross-modal interactions between higher and lower levels of the architecture. MMLatch is motivated by recent advances in cognitive science, analyzing how cognition affects perception and is implemented as a plug and play framework that can be adapted for modern neural architectures. MMLatch improves model performance over our proposed baseline and over MuT. The combination of MMLatch with our baseline achieves state-of-the-art results. We believe top-down cross-modal modeling can augment traditional bottom-up pipelines, improve performance in multimodal tasks and inspire novel multimodal architectures.

In this work, we implement top-down cross-modal modeling as an adaptive feature masking mechanism. In the future, we plan to explore more elaborate implementations that directly affect the state of the network modules from different levels in the network. Furthermore, we aim to extend MMLatch to more tasks, diverse architectures (e.g. Transformers) and for unimodal architectures. Finally, we want to explore the applications top-down masks for model interpretability.

6. REFERENCES

- [1] S. Antol et al., “Vqa: Visual question answering,” in *Proc. CVPR*, 2015.
- [2] Q. You et al., “Image captioning with semantic attention,” in *Proc. CVPR*. IEEE, 2016.
- [3] O. Caglayan, P. Madhyastha, L. Specia, and L. Barrault, “Probing the need for visual context in multimodal machine translation,” in *Proc. NAACL*, 2019.
- [4] G. Paraskevopoulos, S. Parthasarathy, A. Khare, and S. Sundaram, “Multimodal and multiresolution speech recognition with transformers,” in *Proc. 58th ACL*, 2020.

- [5] T. Srinivasan, R. Sanabria, F. Metze, and D. Elliott, "Multimodal speech recognition with unstructured audio masking," in *Proc. 1st Workshop on NLPBT*, 2020.
- [6] J. Klemen and C. D. Chambers, "Current perspectives and methods in studying neural mechanisms of multisensory interactions," *Neuroscience & Biobehavioral Reviews*, 2012.
- [7] P. A. Neil et al., "Development of multisensory spatial integration and perception in humans," *Developmental science*, 2006.
- [8] J. F. Houde and E. F. Chang, "The cortical computations underlying feedback control in vocal production," *Current opinion in neurobiology*, 2015.
- [9] R. L. Shafer et al., "Visual feedback during motor performance is associated with increased complexity and adaptability of motor and neural output," *Behavioural Brain Research*, 2019.
- [10] M. Bar and A. Bubic, "Top-down effects in visual," *The Oxford Handbook of Cognitive Neuroscience, Volume 2*, 2013.
- [11] C. Teufel and B. Nanay, "How to (and how not to) think about top-down influences on visual perception," *Consciousness and Cognition*, 2017.
- [12] E. Sohoglu, J. E. Peelle, R. P. Carlyon, and M. H. Davis, "Predictive top-down integration of prior knowledge during speech perception," *Journal of Neuroscience*, 2012.
- [13] G. Lupyan, "Objective effects of knowledge on visual perception," *Journal of experimental psychology: human perception and performance*, 2017.
- [14] E. Balci et al. and D. Dunning, "Wishful seeing: More desired objects are seen as closer," *Psychological science*, 2010.
- [15] D. R. Proffitt, M. Bhalla, R. Gossweiler, and J. Midgett, "Perceiving geographical slant," *Psychonomic bulletin & review*, 1995.
- [16] S. Manita et al., "A top-down cortical circuit for accurate sensory perception," *Neuron*, 2015.
- [17] C. Firestone and B. J. Scholl, "'top-down' effects where none should be found: The el greco fallacy in perception research," *Psychological science*, 2014.
- [18] C. C. Lee et al., "Emotion recognition using a hierarchical binary decision tree approach," *Speech Communication*, 2011.
- [19] V. Rozgić et al., "Ensemble of svm trees for multimodal emotion recognition," in *Proc. APSIPA. IEEE*, 2012.
- [20] A. Metallinou et al., "Context-sensitive learning for enhanced audiovisual emotion classification," *Transactions on Affective Computing*, 2012.
- [21] M. Wöllmer et al., "Lstm-modeling of continuous emotions in an audiovisual affect recognition framework," *Image and Vision Computing*, 2013.
- [22] A. Shenoy and A. Sardana, "Multilogue-net: A context-aware RNN for multi-modal emotion detection and sentiment analysis in conversation," in *Proc. 2nd Challenge-HML*, 2020.
- [23] S. Poria, I. Chaturvedi, E. Cambria, and A. Hussain, "Convolutional mkl based multimodal emotion recognition and sentiment analysis," in *Proc. ICDM. IEEE*, 2016.
- [24] T. Baltrušaitis, C. Ahuja, and L. P. Morency, "Multimodal machine learning: A survey and taxonomy," *Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [25] A. Zadeh et al., "Tensor fusion network for multimodal sentiment analysis," in *Proc. EMNLP*, 2017.
- [26] Z. Liu et al., "Efficient low-rank multimodal fusion with modality-specific factors," in *Proc. 56th ACL*, 2018.
- [27] A. Zadeh et al., "Multi-attention recurrent network for human communication comprehension," *Proc. AAAI*, 2018.
- [28] A. Zadeh et al., "Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph," in *Proc. 56th ACL*, 2018.
- [29] Y. Gu et al., "Multimodal affective analysis using hierarchical attention strategy with word-level alignment," in *Proc. ACL*, 2018.
- [30] E. Georgiou, C. Papaioannou, and A. Potamianos, "Deep hierarchical fusion with application in sentiment analysis," *Proc. Interspeech*, 2019.
- [31] H. Pham et al., "Found in translation: Learning robust joint representations by cyclic translations between modalities," in *Proc. AAAI*, 2019.
- [32] Z. Sun, P. K. Sarma, W. Sethares, and E. P. Bucy, "Multi-modal sentiment analysis using deep canonical correlation analysis," *Proc. Interspeech*, 2019.
- [33] A. Khare, S. Parthasarathy, and S. Sundaram, "Multi-modal embeddings using multi-task learning for emotion recognition," *Proc. Interspeech*, 2020.
- [34] A. Vaswani et al., "Attention is all you need," in *Proc. 31st NeurIPS*, 2017.
- [35] Y. H. H. Tsai et al., "Multimodal transformer for unaligned multimodal language sequences," in *Proc. 57th ACL*, 2019.
- [36] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Proc. 32nd NeurIPS*, 2019.
- [37] J. B. Delbrouck, N. Tits, M. Brousmiche, and S. Dupont, "A transformer-based joint-encoding for emotion recognition and sentiment analysis," in *2nd Challenge-HML*, 2020.
- [38] W. Rahman et al., "Integrating multimodal information in large pretrained transformers," in *Proc. 58th ACL*, 2020.
- [39] Y. Wang et al., "Words can shift: Dynamically adjusting word representations using nonverbal behaviors," in *Proc. AAAI*, 2019.
- [40] A. Kumar and J. Vepa, "Gated mechanism for attention based multi modal sentiment analysis," in *ICASSP. IEEE*, 2020.
- [41] Y. H. H. Tsai et al., "Multimodal routing: Improving local and global interpretability of multimodal language analysis," in *Proc. EMNLP*, 2020.
- [42] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between capsules," in *Proc. 30th NeurIPS*, 2017.
- [43] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, 1997.
- [44] Diederik P. K. and Jimmy B., "Adam: A method for stochastic optimization," in *3rd ICLR*, Yoshua B. and Yann L., Eds., 2015.
- [45] H. Wen, S. You, and Y. Fu, "Cross-modal context-gated convolution for multi-modal sentiment analysis," *Pattern Recognition Letters*, 2021.
- [46] S. Sourav and J. Ouyang, "Lightweight models for multimodal sequential data," in *Proc. 11th WASSA*, 2021.