

RethinkingTMSC: An Empirical Study for Target-Oriented Multimodal Sentiment Classification

Anonymous ACL submission

Abstract

001 Recently, Target-oriented Multimodal Senti- 042
002 ment Classification (TMSC) has gained signifi- 043
003 cant attention among scholars. However, cur- 044
004 rent multimodal models have reached a perfor- 045
005 mance bottleneck. To investigate the causes of 046
006 this problem, we perform extensive empirical 047
007 evaluation and in-depth analysis of the datasets 048
008 to answer the following questions: **Q1**: Are the 049
009 modalities equally important for TMSC? **Q2**: 050
010 Which multimodal fusion modules are more 051
011 effective? **Q3**: Do existing datasets adequately 052
012 support the research? Our experiments and 053
013 analysis reveal that the current TMSC systems 054
014 primarily rely on the textual modality as most 055
015 of targets' sentiments can be determined *solely* 056
016 by text. Furthermore, many images even lack 057
017 targets in existing datasets. Therefore, con- 058
018 structuring a more suitable dataset for TMSC is 059
019 urgently needed since it has seriously hindered 060
020 the research progress. 061

021 1 Introduction 062

022 Target-oriented sentiment classification, also 063
023 known as aspect-based sentiment classification, is 064
024 a fundamental task of sentiment analysis (Pontiki 065
025 et al., 2014, 2015, 2016). It aims to judge the sen- 066
026 timental polarity (positive, negative, or neutral) of 067
027 a specific target within text. To improve the per- 068
028 formance by considering multimodal information, 069
029 Target-oriented Multimodal Sentiment Classifica- 070
030 tion (TMSC) is proposed to integrate both visual 071
031 and textual information (Yu and Jiang, 2019). The 072
032 major challenge of this task is how to utilize infor- 073
033 mation from multiple modalities effectively. 074

034 Recently, the performance of the TMSC systems 075
035 gradually reaches a plateau and the progress in 076
036 tackling this task has slowed down. Using the F1- 077
037 score metric on the popular datasets, Twitter15 and 078
038 Twitter17 (Yu and Jiang, 2019), we observe that 079
039 state-of-the-art baselines only achieve an F1-score 080
040 of around 70. Therefore, in this paper, we aim to 081
041 analyze the causes behind it at both model level 082

and modality level. Roughly speaking, the modules 042
in the model structures can be categorized into two 043
types: 1) encoders to model the representations 044
of different modalities; and 2) multimodal fusion 045
modules to model the interactions between modal- 046
ities. Moreover, we give a deep analysis of the 047
characteristics of two widely-used datasets, aiming 048
to answer the following three questions: 049

Q1: Are the modalities equally important for 050
TMSC? To explore this issue, we compare and 051
analyze the performance of unimodal models on 052
this task. For the textual modality, we use BERT 053
(Devlin et al., 2019) as the backbone, as it is a 054
widely-used pre-trained language model outper- 055
forming earlier models like LSTM (Hochreiter and 056
Schmidhuber, 1997), memory network (Weston 057
et al., 2015), etc. (Yu and Jiang, 2019). For the vi- 058
sual modality, ResNet (He et al., 2016), ViT (Doso- 059
vitskiy et al., 2021), and Faster R-CNN (Ren et al., 060
2015) are adopted to learn powerful representations 061
(see Figure 1). 062

Q2: Which multimodal fusion modules are more 063
effective? The current models use various fusion 064
strategies to model the interactions between modal- 065
ities, while obtaining little improvement (less than 066
2 points of F1 score) (see Table 1). To explore 067
the effectiveness of different fusion approaches, we 068
summarize the fusion strategies into six categories: 069
Concatenation, Tensor Fusion (Zadeh et al., 2017), 070
Self Attention (Vaswani et al., 2017), Image2Text, 071
Text2Image and Bi-direction. Then we perform a 072
comparative study of these fusion modules using a 073
unified setup to eliminate potential bias from model 074
size and structure (see Figure 2). 075

Q3: Do existing datasets adequately support the 076
research? We analyze the existing datasets (i.e., 077
Twitter15 and Twitter17) in depth to explore their 078
limitations and obtain the following findings: 1) 079
The size of existing datasets is limited and the dis- 080
tribution of the sentiments is unbalanced; 2) The 081
multimodal sentiment is much more consistent with 082

Model	Image Encoder			Fusion Module			Performance (F1)			
	ResNet	ViT	Faster R-CNN	Concat	Tensor Fusion	Self Attention	Image2Text	Text2Image	Twitter15	Twitter17
Res-BERT+BL	✓	✗	✗	✓	✗	✓	✗	✗	69.21	66.48
Res-BERT+BL-TFN	✓	✗	✗	✗	✓	✓	✗	✗	68.74	64.29
mBERT	✓	✗	✗	✓	✗	✓	✗	✗	71.07	67.06
TomBERT	✓	✗	✗	✓	✗	✓	✓	✗	71.75	68.03
EF-CapTrBERT	✓	✗	✗	✗	✗	✓	✗	✗	73.25	68.42
SMP	✗	✓	✗	✗	✗	✗	✓	✓	72.24	69.47
VLP	✗	✗	✓	✗	✗	✓	✗	✗	73.80	71.80

Table 1: The model structures of various baselines.

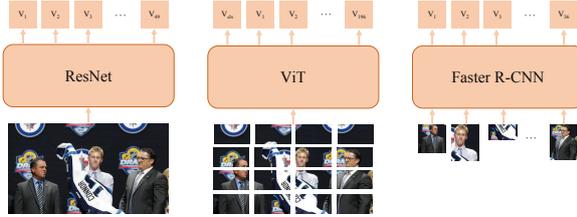


Figure 1: Different image encoders.

the textual sentiment than the visual sentiment; 3) A large number of targets do not exist in images; 4) There are only a small number of samples where the sentiment is decided by both text and image.

The main contributions of this work are as follows: 1) We investigate the effectiveness of different model structures for TMSC, including various unimodal encoders and multimodal fusion modules; 2) We give an in-depth analysis of limitations of existing widely-used datasets; 3) We derive several valuable observations and point out promising directions for the future research of TMSC model design and dataset creation.

2 Empirical Study

We summarize the model structures and performance of the baselines for the TMSC task in Table 1. Their structural differences are mainly reflected in the different unimodal encoders and multimodal fusion modules used. Therefore, we carry out several experiments to analyze the impact of these two aspects on performance.

2.1 Unimodal Encoder

As previously mentioned in Section 1, we primarily focus on exploring the different image encoders, ResNet, ViT, and Faster R-CNN (see Figure 1), while using BERT as the text encoder.

ResNet. Following most of the baselines (e.g., mBERT (Yu and Jiang, 2019), TomBERT (Yu and Jiang, 2019) and EF-CapTrBERT (Khan and Fu, 2021)), we adopt ResNet-152 as one of the image encoders. Each image is resized into 224 by 224,

and then passed through the model to obtain 49 regions, which are used as the image representation $I = [v_1, v_2, \dots, v_{49}]$, where $v_i \in \mathbb{R}^{2048}$.

ViT. Being same as SMP (Ye et al., 2022), we adopt ViT to model the image by dividing it into 16 by 16 patches. A CLS token is added at the beginning and fed into the Transformer (Vaswani et al., 2017) encoder to obtain the image representation $I = [v_{cls}, v_1, v_2, \dots, v_{196}]$, where $v_i \in \mathbb{R}^{768}$.

Faster R-CNN. Similar to VLP (Ling et al., 2022), we adopt Faster R-CNN that is retrained on the Visual Genome Dataset (Krishna et al., 2017). We select the top 36 object proposals as the image representation $I = [v_1, v_2, \dots, v_{36}]$, where $v_i \in \mathbb{R}^{2048}$ is obtained from the ROI pooling layer of the Region Proposal Network (Ren et al., 2015).

2.2 Multimodal Fusion

We categorize the current multimodal fusion modules into five groups as follows (see Figure 2).

Concatenate. Concatenate is the simplest form of fusion, where the pooled text representation $H_p^T \in \mathbb{R}^{768}$ is directly combined with the pooled image representation $H_p^I \in \mathbb{R}^{768}$ ¹ to obtain the multimodal representation $H = H_p^I \oplus H_p^T$, where \oplus is a concatenation operation and $H \in \mathbb{R}^{768+768}$.

Tensor Fusion. Tensor Fusion is proposed for modeling interactions between modalities while preserving the characteristics of individual modalities. Through this module, we obtain the multimodal representation $H = H_p^I \otimes H_p^T$, where \otimes is an outer product operation and $H \in \mathbb{R}^{768 \times 768}$.

Self Attention. Self Attention first concatenates the image representation $H^I \in \mathbb{R}^{l_I \times 768}$ with the text representation $H^T \in \mathbb{R}^{l_T \times 768}$, where l_I and l_T are the lengths of image and text, respectively. Then the concatenated representation $H^I \oplus H^T$ is passed through three self-attention layers and a pooling layer to obtain the multimodal representa-

¹A linear mapping layer is added after the image encoder to map the image representation to 768 dimensions to ensure uniformity when using different image encoders.

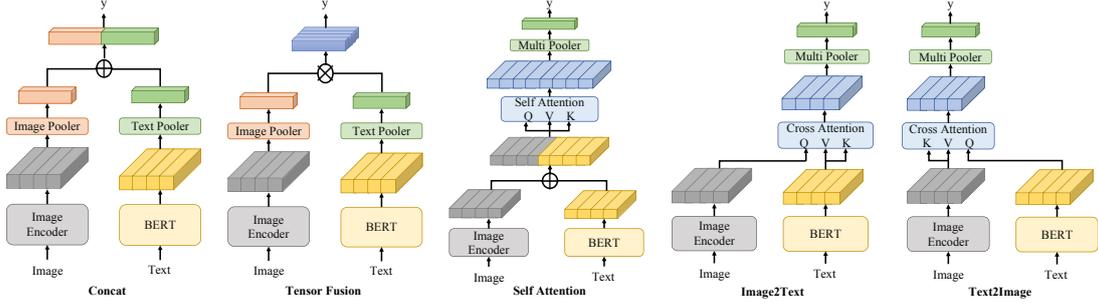


Figure 2: Various multimodal fusion modules.

Modality	Model	Twitter15		Twitter17	
		ACC	F1	ACC	F1
Text	BERT	76.72±1.16	71.19±2.19	68.04±0.40	65.66±0.35
Image	ResNet	57.65±1.00	32.52±2.66	57.79±0.99	51.98±1.23
	ViT	59.65±1.13	31.25±2.71	59.53±0.95	54.08±0.78
	Faster R-CNN	55.97±1.10	<u>35.72±5.43</u>	56.18±0.85	49.88±1.70
	ResNet				
	Concatenate	75.29±0.45	68.71±1.34	67.92±0.56	65.32±0.53
	Tensor Fusion	74.19±0.94	68.93±0.57	66.66±1.21	63.99±1.61
	Self Attention	76.03±0.96	70.57±2.39	68.01±0.96	65.41±1.60
	Image2Text	77.13±1.33	71.48±1.90	<u>69.37±0.36</u>	<u>66.85±0.79</u>
	Text2Image	75.18±1.66	67.77±4.81	68.07±0.58	65.18±1.48
	Bi-direction	<u>77.32±0.63</u>	72.06±0.81	68.41±1.01	66.39±1.39
	ViT				
Multimodal	Concatenate	76.22±0.90	70.37±1.45	67.94±0.70	66.17±0.78
	Tensor Fusion	73.44±0.78	67.46±1.45	65.46±1.67	62.02±1.40
	Self Attention	75.08±0.41	68.94±0.83	67.52±0.58	65.56±0.35
	Image2Text	<u>77.11±0.44</u>	<u>71.91±0.42</u>	69.14±0.52	66.96±0.68
	Text2Image	75.12±1.01	69.40±1.38	67.52±1.06	64.49±1.46
	Bi-direction	76.70±0.75	71.67±1.45	<u>69.16±0.17</u>	<u>67.25±0.56</u>
	Faster R-CNN				
	Concatenate	75.45±0.73	69.77±1.23	67.60±1.15	64.74±1.69
	Tensor Fusion	72.09±0.66	66.77±1.04	66.34±1.45	62.96±2.09
	Self Attention	76.09±0.89	70.08±1.37	68.09±1.10	66.12±1.23
	Image2Text	77.36±0.37	<u>71.69±0.37</u>	68.43±0.65	66.44±1.10
	Text2Image	70.82±2.99	57.94±5.81	60.31±6.43	54.50±7.06
	Bi-direction	76.57±0.46	70.88±0.89	69.51±0.62	67.50±0.37

Table 2: TMSC results on Twitter15 and Twitter17 datasets. The overall best results and those within each corresponding block are marked with **bold** and underline, respectively.

tion $H \in \mathbb{R}^{768}$.

Image2Text. Image2Text is a kind of cross-attention mechanism (Vaswani et al., 2017), using the image representation H^I as the query and the text representation H^T as the key and value, through three attention layers to get the multimodal representation $H \in \mathbb{R}^{768}$.

Text2Image. The only difference between Text2Image and Image2Text is that Text2Image uses the text representation H^T as the query and the image representation H^I as the key and value.

Bi-direction. Bi-direction means that the representations of Text2Image and Image2Text are concatenated as the multimodal representation $H \in \mathbb{R}^{768+768}$.

2.3 Results Analysis

We perform experiments of different unimodal encoders and fusion modules over Twitter15 and Twitter17. For each set of experiments, we test under five random seeds (i.e., 0, 42, 199, 2022, and 11122) and present the mean and standard deviation as the final result. The Adam optimizer (Kingma and Ba, 2015) is used with a learning rate of $2e-5$ and each experiment is run on a 3090 GPU for 8 epochs. In Table 2, we show the results and we have the following observations:

First, the text-only model (i.e., BERT) performs well, while the visual-only models (i.e., ResNet, ViT, and Faster R-CNN) perform relatively poorly, revealing that the reliance on text is much greater than that on images for the TMSC task on these two datasets. In comparison, this phenomenon is more pronounced in Twitter15.

Second, the performance of the model is affected by the use of different fusion methods. Specifically, fusion modules that primarily focus on acquiring the textual information (e.g., Image2Text) perform better than those focused on acquiring the visual information (e.g., Text2Image). This again reveals the inconsistent importance of text and images.

Third, compared with the text-only model, the various fusion modules do not have significant gains in performance and some are even worse. This is due to the fact that some images do not provide related information, but rather distracting information instead. We give a detailed analysis of the performance comparison for the multimodal model versus the text-only model in Appendix B.

Fourth, the impact of various image encoders is not clear, as evidenced by low performance and high standard deviation on the two datasets (see the ‘‘Image’’ part of Table 2). Moreover, differences in performance among various image encoders are small in the multimodal fusion settings (see the ‘‘Multimodal’’ part of Table 2). This is due to the

Dataset	Twitter15					Twitter17				
	#Negative	#Neutral	#Positive	#Total	#Avg Targets	#Negative	#Neutral	#Positive	#Total	#Avg Targets
Train	368	1883	928	3179	1.348	416	1638	1508	3562	1.410
Dev	149	670	303	1122	1.336	144	517	515	1176	1.439
Test	113	607	317	1037	1.354	168	573	493	1234	1.450

Table 3: Statistics of the datasets. #Avg Targets means the average number of targets for each sample.

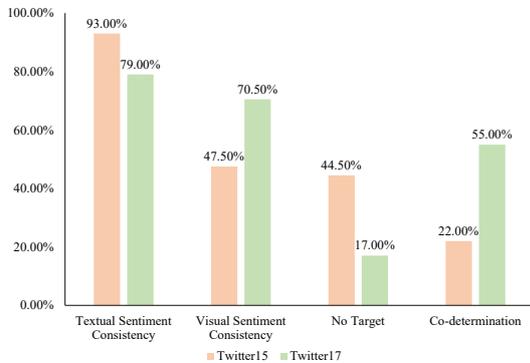


Figure 3: Annotation analysis. Textual/Visual Sentiment Consistency: the consistency of the target’s sentiment in text/image with the sentiment in multimodal information. No Target: the percentage of images that are missing the target for sentiment analysis. Co-determination: the percentage of targets that sentiment is jointly determined by text and image.

characteristics of visual data in existing datasets, which is analyzed in depth in the following section.

3 Data Analysis

To gain a deeper understanding of the performance issues mentioned above, we analyze the two datasets in depth. Following the annotation procedure by Yu and Jiang (2019), we invite three domain experts to annotate 400 samples randomly selected from the test set (200 from Twitter15 and 200 from Twitter17) with four aspects and take the majority vote as the final annotation result (Figure 3)². We have the following observations:

First, as shown in Table 3, the sample size is relatively small, with an average of less than 1.5 targets per sample. Additionally, the distributions of the sentimental labels are unbalanced in both datasets, with neutral sentiment accounting for approximately 50% and negative sentiment accounting for less than 15%. The reason behind this is that Twitter15 and Twitter17 were originally constructed by Zhang et al. (2018) and Lu et al. (2018) respectively for the named entity recognition task, rather than specifically for TMSC.

Second, the multimodal sentiment has high consistency with the textual sentiment but low consistency with the visual sentiment. In Twitter15, 93%

²Illustrative examples with annotations are in Appendix C.

of the targets have the same textual sentiment as the multimodal sentiment, while only 47.5% have a visual sentiment that matches. This indicates the biased distribution existing in the dataset, i.e., the textual information is more useful for determining the multimodal sentiment. Although this phenomenon is mitigated in Twitter17, the textual information is still more consistent with the multimodal sentiment than the visual information.

Third, a large number of targets do not exist in images, which is also not suitable for the *target-oriented* multimodal sentiment classification task. This problem may stem from the construction of the two datasets, where the targets are selected directly from the text, without taking into account the corresponding images (Yu and Jiang, 2019).

Fourth, due to the problems of irrelevant images and non-existence of targets in images, there is only a small portion of the data where the multimodal sentiment is determined by both text and images. Specifically, only 22% of Twitter15 and 55% of Twitter17 data require both text and images for the sentiment classification. As for the multimodal task, these two datasets may not be the best-suited in this aspect.

Taken together, text plays a dominant role in the existing TMSC benchmark, which is consistent with the results analyzed in Section 2. In addition, it also has problems such as unbalanced label distribution and missing targets, which has seriously hindered TMSC research.

4 Conclusion and Future Work

In this paper, we conduct a series of in-depth experiments for TMSC and data analysis of existing datasets. We find that the current multimodal models do not achieve significant performance gains compared to the text-only models, primarily due to the fact that existing datasets rely heavily on the textual modality while the visual modality seems to be less important. In particular, existing datasets suffer from missing targets in images, which cannot fully support TMSC studies. Therefore, it is crucial to build more suitable datasets to promote the progress of tackling this task.

277 Limitations

278 Although we have conducted a series of experi-
279 ments and data analysis for the TMSO task to the
280 best of our ability, there are at least the following
281 limitations to our work. First, our data analysis
282 was performed mainly for the currently publicly
283 available English datasets Twitter15 and Twitter17,
284 neglecting the Chinese dataset Multi-ZOL, which
285 has not been widely studied. Second, although our
286 analysis indicated some problems in using the cur-
287 rently dataset to measure the TMSO task, we did
288 not construct a new and better dataset for use in
289 academic studies. We have included this task as
290 one of our future works to be investigated.

291 References

292 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
293 Kristina Toutanova. 2019. [BERT: pre-training of
294 deep bidirectional transformers for language under-
295 standing](#). In *Proceedings of the 2019 Conference of
296 the North American Chapter of the Association for
297 Computational Linguistics: Human Language Techno-
298 logies, NAACL-HLT 2019, Minneapolis, MN, USA,
299 June 2-7, 2019, Volume 1 (Long and Short Papers)*,
300 pages 4171–4186. Association for Computational
301 Linguistics.

302 Alexey Dosovitskiy, Lucas Beyer, Alexander
303 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
304 Thomas Unterthiner, Mostafa Dehghani, Matthias
305 Minderer, Georg Heigold, Sylvain Gelly, Jakob
306 Uszkoreit, and Neil Houlsby. 2021. [An image
307 is worth 16x16 words: Transformers for image
308 recognition at scale](#). In *9th International Conference
309 on Learning Representations, ICLR 2021, Virtual
310 Event, Austria, May 3-7, 2021*. OpenReview.net.

311 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian
312 Sun. 2016. [Deep residual learning for image recogni-
313 tion](#). In *2016 IEEE Conference on Computer Vision
314 and Pattern Recognition, CVPR 2016, Las Vegas,
315 NV, USA, June 27-30, 2016*, pages 770–778. IEEE
316 Computer Society.

317 Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long
318 short-term memory](#). *Neural Comput.*, 9(8):1735–
319 1780.

320 Zaid Khan and Yun Fu. 2021. [Exploiting BERT for
321 multimodal target sentiment classification through
322 input space translation](#). In *MM '21: ACM Multimedia
323 Conference, Virtual Event, China, October 20 - 24,
324 2021*, pages 3034–3042. ACM.

325 Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A
326 method for stochastic optimization](#). In *3rd Inter-
327 national Conference on Learning Representations,
328 ICLR 2015, San Diego, CA, USA, May 7-9, 2015,
329 Conference Track Proceedings*.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin John- 330
son, Kenji Hata, Joshua Kravitz, Stephanie Chen, 331
Yannis Kalantidis, Li-Jia Li, David A. Shamma, 332
Michael S. Bernstein, and Li Fei-Fei. 2017. [Vi- 333
sual genome: Connecting language and vision us-
334 ing crowdsourced dense image annotations](#). *Int. J.
335 Comput. Vis.*, 123(1):32–73. 336

Yan Ling, Jianfei Yu, and Rui Xia. 2022. [Vision- 337
language pre-training for multimodal aspect-based
338 sentiment analysis](#). In *Proceedings of the 60th An-
339 nual Meeting of the Association for Computational
340 Linguistics (Volume 1: Long Papers), ACL 2022,
341 Dublin, Ireland, May 22-27, 2022*, pages 2149–2159.
342 Association for Computational Linguistics. 343

Di Lu, Leonardo Neves, Vitor Carvalho, Ning Zhang, 344
and Heng Ji. 2018. [Visual attention model for name
345 tagging in multimodal social media](#). In *Proceedings
346 of the 56th Annual Meeting of the Association for
347 Computational Linguistics, ACL 2018, Melbourne,
348 Australia, July 15-20, 2018, Volume 1: Long Papers*,
349 pages 1990–1999. Association for Computational
350 Linguistics. 351

Maria Pontiki, Dimitris Galanis, Haris Papageorgiou, 352
Ion Androutsopoulos, Suresh Manandhar, Moham- 353
mad Al-Smadi, Mahmoud Al-Ayyoub, Yanyan Zhao,
354 Bing Qin, Orphée De Clercq, Véronique Hoste,
355 Marianna Apidianaki, Xavier Tannier, Natalia V.
356 Loukachevitch, Evgeniy V. Kotelnikov, Núria Bel,
357 Salud María Jiménez Zafra, and Gülsen Eryigit. 2016.
358 [Semeval-2016 task 5: Aspect based sentiment analy-
359 sis](#). In *Proceedings of the 10th International Work-
360 shop on Semantic Evaluation, SemEval@NAACL-
361 HLT 2016, San Diego, CA, USA, June 16-17, 2016*,
362 pages 19–30. The Association for Computer Linguis-
363 tics. 364

Maria Pontiki, Dimitris Galanis, Haris Papageorgiou, 365
Suresh Manandhar, and Ion Androutsopoulos. 2015.
366 [Semeval-2015 task 12: Aspect based sentiment analy-
367 sis](#). In *Proceedings of the 9th International Work-
368 shop on Semantic Evaluation, SemEval@NAACL-
369 HLT 2015, Denver, Colorado, USA, June 4-5, 2015*,
370 pages 486–495. The Association for Computer Lin-
371 guistics. 372

Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Har- 373
ris Papageorgiou, Ion Androutsopoulos, and Suresh 374
Manandhar. 2014. [Semeval-2014 task 4: Aspect
375 based sentiment analysis](#). In *Proceedings of the 8th
376 International Workshop on Semantic Evaluation, Sem-
377 Eval@COLING 2014, Dublin, Ireland, August 23-
378 24, 2014*, pages 27–35. The Association for Com-
379 puter Linguistics. 380

Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian 381
Sun. 2015. [Faster R-CNN: towards real-time ob-
382 ject detection with region proposal networks](#). In *Ad-
383 vances in Neural Information Processing Systems 28:
384 Annual Conference on Neural Information Process-
385 ing Systems 2015, December 7-12, 2015, Montreal,
386 Quebec, Canada*, pages 91–99. 387

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Jason Weston, Sumit Chopra, and Antoine Bordes. 2015. [Memory networks](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Nan Xu, Wenji Mao, and Guandan Chen. 2019. [Multi-interactive memory network for aspect based multimodal sentiment analysis](#). In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 371–378. AAAI Press.

Junjie Ye, Jie Zhou, Junfeng Tian, Rui Wang, Jingyi Zhou, Tao Gui, Qi Zhang, and Xuanjing Huang. 2022. [Sentiment-aware multimodal pre-training for multimodal sentiment analysis](#). *Knowl. Based Syst.*, 258:110021.

Jianfei Yu and Jing Jiang. 2019. [Adapting BERT for target-oriented multimodal sentiment classification](#). In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 5408–5414. ijcai.org.

Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. [Tensor fusion network for multimodal sentiment analysis](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 1103–1114. Association for Computational Linguistics.

Qi Zhang, Jinlan Fu, Xiaoyu Liu, and Xuanjing Huang. 2018. [Adaptive co-attention network for named entity recognition in tweets](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5674–5681. AAAI Press.

A Related Work

As one of the tasks of sentiment analysis, TMSC has gained great attention from scholars in recent years (Yu and Jiang, 2019). Xu et al. (2019) constructed a Chinese dataset named Multi-ZOL and proposed a multi-hop memory network for handling modal interactions. Subsequently, Yu and Jiang (2019) constructed two English datasets, Twitter15 and Twitter17, and applied BERT to this

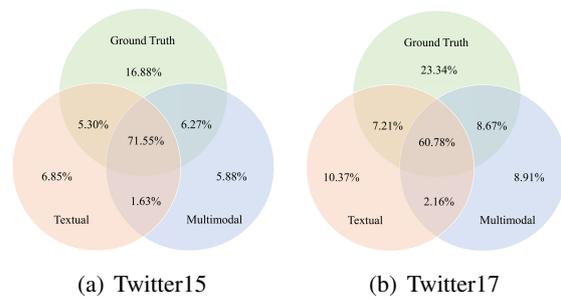


Figure 4: Venn diagram for model performance visualization.

task. The following research on the TMSC task can be divided into two directions. On the one hand, there is the continuous exploration of how to enhance the interactions between modalities (Khan and Fu, 2021), and on the other hand, there is the application of pre-trained models to this task (Ye et al., 2022; Ling et al., 2022). Despite these efforts, the current models have not yet achieved significant performance gains relative to the text-only models. We have conducted a series of experiments and data analysis, hoping to provide some insights for the future research of TMSC.

B Model Performance Visualization

We select Image2Text (Faster R-CNN) as the representative of the multimodal models and compare its performance with that of BERT with a random seed of 11122 to obtain Figure 4. The intersection of every two circles in the figure represents the part where the prediction results are consistent. Based on this comparison, we have the following observations: **First**, in terms of prediction accuracy, the multimodal model does not achieve a significant gain over the text-only model. **Second**, a portion of the data is predicted correctly by the multimodal model but incorrectly by the text-only model, and vice versa. The proportions of these two parts are similar. This suggests that when images do contribute valuable information to the multimodal model, they also introduce noise. In order to improve the performance, further investigation is required for how to properly incorporate the visual information. **Third**, over 16% of the data has sentiments that neither the text-only model nor the multimodal model predicts correctly. This indicates the weakness of the current models and we need further explorations.

Text	Image	Target	Sentiment		
			Multimodal	Textual	Visual
Congratulations to our second draw winner - Bulaire Leber of ADSS Global , Haiti . Thanks for participating , Bulaire		Bulaire Leber	Positive	Positive ✓	Positive ✓
RT @ BeschlossDC : Coretta Scott King with Robert amp Ethel Kennedy after husband ' s assassination , which occurred tonight 1968		Ethel Kennedy	Negative	Negative ✓	Neutral ✗
Pres Obama takes the stage at @ RutgersU Commencement in school football stadium in Piscataway , NJ .		Obama	Positive	Neutral ✗	Positive ✓
RT @ Refugees : Today , 18 - year - old Yehya became the 1 millionth Syrian to register as a refugee in Lebanon		Lebanon	Neutral	Neutral ✓	<u>No Target</u> ✗

Table 4: The annotation examples.

C Annotation Examples

To clearly and visually illustrate the various scenarios that arise during the dataset annotation process, four samples are presented in Table 4.

The **first** example demonstrates a scenario where the textual sentiment and the visual sentiment matches, resulting in a multimodal sentiment determined by both modalities. In the example, the sentiment in the text is determined to be positive through the use of words such as “Congratulations” and “winner”. Similarly, the sentiment in the image can be inferred as positive by identifying the target (i.e., the first person on the right) and noticing his smiling face.

The **second** example shows a scenario where the textual sentiment aligns with the multimodal sentiment but not with the visual sentiment, leading to a multimodal sentiment determined by the textual modality only. Specifically, the sentiment conveyed by the text is negative due to the phrase

“after husband’s assassination” and the sentiment conveyed by the image is neutral as it does not show an obvious facial expression on the person referred to in the text (i.e., the first person on the left). Therefore, the multimodal sentiment conveyed by both modalities is negative.

Corresponding to the second example, the **third** example illustrates a scenario where the visual sentiment aligns with the multimodal sentiment but not with the textual sentiment. In particular, the text simply states a fact with a neutral sentiment, while the image shows the target (i.e., the person waving his hand in front of the podium) with a positive facial expression and posture, resulting in a positive multimodal sentiment overall.

The **fourth** example presents a scenario where there is no target in the image, resulting in a multimodal sentiment determined solely by the textual modality. Here, the target is “Lebanon”, but since there is only one person in the image and no information about “Lebanon”, we can only conclude

520 that the multimodal sentiment is neutral based on
521 the text. It is worth mentioning that such a sample
522 is not ideal for the TMSA task as the image does
523 not convey any sentimental information towards
524 the target.