

# Aspect-specific Context Modeling for Aspect-based Sentiment Analysis

Fang Ma, Chen Zhang, Bo Zhang, Dawei Song\*

School of Computer Science, Beijing Institute of Technology  
Beijing, China

{mfang, czhang, bo.zhang, dwsong}@bit.edu.cn

## Abstract

Aspect-based sentiment analysis (ABSA) aims at predicting sentiment polarity (SC) or extracting opinion span (OE) expressed towards a given aspect. Previous work in ABSA mostly relies on rather complicated aspect-specific feature induction. Recently, pretrained language models (PLMs), e.g., BERT, have been used as context modeling layers to simplify the feature induction structures and achieve state-of-the-art performance. However, such PLM-based context modeling can be not that aspect-specific. Therefore, a key question is left under-explored: how the aspect-specific context can be better modeled through PLMs? To answer the question, we attempt to enhance aspect-specific context modeling with PLM in a non-intrusive manner. We propose three aspect-specific input transformations, namely aspect companion, aspect prompt, and aspect marker. Informed by these transformations, non-intrusive aspect-specific PLMs can be achieved to promote the PLM to pay more attention to the aspect-specific context in a sentence. Additionally, we craft an adversarial benchmark for ABSA (advABSA) to see how aspect-specific modeling can impact model robustness. Extensive experimental results on standard and adversarial benchmarks for SC and OE demonstrate the effectiveness and robustness of the proposed method, yielding new state-of-the-art performance on OE and competitive performance on SC.<sup>1</sup>

## 1 Introduction

Aspect-based sentiment analysis (ABSA) aims to infer multiple fine-grained sentiments from the same content, with respect to multiple aspects. A fine-grained sentiment in ABSA can be categorized into two forms, i.e., sentiment and opinion. Accordingly, two sub-tasks of ABSA are aspect-based

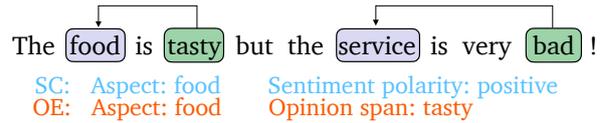


Figure 1: Example of the SC and OE. The words highlighted in purple represent the given aspects, whereas the words in green represent the corresponding opinion.

sentiment classification (SC for short) and aspect-based opinion extraction (OE for short). Given an aspect in a sentence, SC aims to predict its sentiment polarity, while OE aims to extract the corresponding opinion span expressed towards the given aspect. Figure 1 shows an example of SC and OE. In the sentence “*The food is tasty but the service is very bad!*”, if *food* is the given aspect, SC requires a model to give a *positive* sentiment on *food* while OE requires a model to extract *tasty* as the opinion span for the aspect *food*.

An effective ABSA model typically would require either aspect-specific feature induction or context modeling. Prior work in ABSA largely relies on rather complicated aspect-specific feature induction to achieve a good performance. Recently, pretrained language models (PLMs) have been shown to enhance the state-of-the-art ABSA models due to their extraordinary context modeling ability. However, currently the use of PLMs in these ABSA models is aspect-general, but overlooks two key questions: 1) whether the context modeling of a PLM can be aspect-specific; and 2) whether the aspect-specific context modeling within a PLM can further enhance ABSA.

To address the aforementioned key questions, in this paper, we propose to achieve *aspect-specific context modeling* of PLMs with *aspect-specific input transformations*. In addition to the commonly used aspect-specific input transformation that appends an aspect to a sentence, i.e., **aspect companion**, we propose two more aspect-specific in-

\*Dawei Song is the corresponding author.

<sup>1</sup>The code and proposed data are available at <https://github.com/BD-MF/ASCM4ABSA>.

put transformations, namely **aspect prompt** and **aspect marker**, to explicitly mark a concerned aspect in a sentence. Aspect prompt shares a similar idea with aspect companion, except that it appends an aspect-oriented prompt instead of sole aspect description to the sentence. Aspect marker distinguishes itself from the above two by introducing two marker tokens, one before and the other after the aspect. As the proposed input transformations are intended to highlight a specific aspect, they in turn can be leveraged to promote the PLM to pay more attention to the context that is relevant to the aspect. Methodologically, this is achieved with a novel aspect-focused PLM fine-tuning model that is guided by the input transformations and essentially performs a joint context modeling and aspect-specific feature induction.

We conduct extensive experiments on both sub-tasks of ABSA, i.e., SC and OE, with various standard benchmarking datasets for effectiveness test, along with our crafted adversarial ones for robustness test. Since there are only datasets for robustness tests in SC and is currently no dataset for robustness tests in OE, we propose an adversarial benchmark (advABSA) based on (Xing et al., 2020)’s datasets and methods. That is, the advABSA benchmark can be decomposed to two parts, where the first part is ARTS-SC for SC reused from (Xing et al., 2020) and the second part is ARTS-OE for OE crafted by us. The results show that models with aspect-specific context modeling achieve the state-of-the-art performance on OE and also outperform various strong SC baseline models without aspect-specific modeling. Overall, these results indicate that aspect-specific context modeling for PLMs can further enhance the performance of ABSA.

To better understand the effectiveness of the three input transformations, we carry out a series of further analyses. After injecting aspect-specific input transformations into a sentence, we observe that the model attends to the correct opinion spans. Hence, we expect that a simple model with aspect-specific context modeling yet without needing complicated aspect-specific feature induction would serve as a sufficiently strong approach for ABSA.

## 2 Related Work

### 2.1 Aspect-based Sentiment Classification SC

ABSA falls in the broad scope of fine-grained opinion mining. As a sub-task of ABSA, SC deter-

mines the sentiment polarity of a given aspect in a sentence and has recently emerged as an active research area with lots of aspect-specific feature induction approaches. These approaches range from memory networks (Tang et al., 2016; Wang et al., 2018), convolutional networks (Li et al., 2018; Huang and Carley, 2018), attentional networks (Wang et al., 2016; Ma et al., 2017), to graph-based networks (Zhang et al., 2019a,b; Wang et al., 2020; Tang et al., 2020). More recently, PLMs such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019), have been applied to SC in a context-encoder scheme (YU and JIANG, 2019; Li et al., 2019; Xu et al., 2019; Liang et al., 2019; Song et al., 2020; Yadav et al., 2021) and achieved the state-of-the-art performance. However, PLMs in these models are aspect-general. We aim to achieve aspect-specific context modeling with PLMs so that these models can be further improved.

### 2.2 Aspect-based Opinion Extraction OE

OE is another sub-task of ABSA, first proposed by Fan et al. (2019). It aims to extract from a sentence the corresponding opinion span describing an aspect. Most work in this area treats OE as a sequence tagging task, for which complex methods are developed to capture the interaction between the aspect and the context (Fan et al., 2019; Wu et al., 2020; Feng et al., 2021). More recent models such as TSMSA-BERT (Feng et al., 2021) and ARGCN-BERT (Jiang et al., 2021), adopt PLMs. In TSMSA-BERT, the multi-head self-attention is utilized to enhance the BERT PLM. ARGCN-BERT uses an attention-based relational graph convolutional network with BERT to exploit syntactic information. We will incorporate our aspect-specific context modeling methods into PLMs to see whether the proposed methods can further improve the OE performance.

## 3 Aspect-specific Context Modeling

### 3.1 Task Description

ABSA (Both SC and OE) requires a pre-given aspect. Formally, a sentence is depicted as  $S = \{w_1, w_2, \dots, w_n\}$  that contains  $n$  words including the aspect. The aspect  $A = \{a_1, a_2, \dots, a_m\}$  is composed of  $m$  words. The goal of SC is to find the sentiment polarity with respect to the given aspect  $A$ . OE aims to extract corresponding opinion span based on the given aspect  $A$ . Recap the example in Figure 1 that contains aspect *food*. SC



aspect, is,  $a_1, \dots, a_m, [\text{SEP}]$ . This format sequence prompts the PLM to target at the aimed aspect.

### 3.4.3 Aspect Marker

Aspect marker inserts markers into the sentence to explicitly mark the boundaries of the concerned aspect. Specifically, we define the markers as two preserved tokens:  $\langle \text{asp} \rangle$  and  $\langle / \text{asp} \rangle$ . We insert them into the input sentence before and after the concerned aspect, to mark the start and end of the given aspect.  $\langle \text{asp} \rangle$  indicates the start of the aspect, and  $\langle / \text{asp} \rangle$  indicates the end of the aspect. Let  $\hat{S}$  denote the modified sequence with aspect marker inserted:  $\hat{S} = \{[\text{CLS}], w_1, \dots, \langle \text{asp} \rangle, a_1, \dots, a_m, \langle / \text{asp} \rangle, \dots, w_n, [\text{SEP}]\}$ .

The three *aspect-specific input transformations* gain significant improvement in our experiments (Section 5), and this strengthens our hypothesis that injecting the aspect marker at the input layer can help the PLM capture aspect-specific contextual information further.

## 3.5 Context Modeling

Previous PLM-based ABSA work directly adopts the hidden states of the PLM for downstream classification. However, an empirical observation is that the context words close to the aspect are more semantic-relevant to the aspect (Ma et al., 2021). In the case, more sentiment information is possibly contained in the aspect’s local context rather than the global context. As a result, the general usage of the hidden states from the PLM loses much local contextual information related to the aspect. With the help of the three input transformations, we obtain the hidden states that incorporate the aspect-oriented local context. Let

$$H = \text{PLM}(\hat{S}) \quad (1)$$

where  $H = \{h_1, h_2, \dots, h_1^a, \dots, h_m^a, \dots, h_n\}$  represents the sequence of hidden states.

## 3.6 Feature Induction

As aforementioned, aspect-general feature induction contains the semantic information critical to the whole sentence rather than the given aspect, and the induced aspect-general feature may be aspect-irrelevant when the sentence contains two or more aspects. After getting the global contextual representation  $H$ , existing work needs an aspect-specific

feature extraction strategy to induce the aspect feature after getting the global contextual representation  $H$ . For an enriched aspect-awareness, we adopt the mean pool on the hidden states corresponding to the first and last aspect tokens. Let

$$\hat{H} = \text{MeanPool}([h_1^a, h_m^a]) \quad (2)$$

represent the aspect-specific feature, where  $h_1^a$  indicates the hidden state of the first aspect token, and  $h_m^a$  indicates the hidden state of the last aspect token. Due to that OE is a token-level classification task, we concatenate the aspect-specific feature  $\hat{H}$  and the global contextual representation  $H$  as the final aspect-specific contextual representation for tagging.

## 3.7 Fine-tuning

After getting the aspect-specific contextual representation  $\hat{H}$ , an multi-layered Perceptron (MLP) layer is used to fine-tune the proposed BERT or RoBERTa based model. The MLP contains four steps: a fully-connected layer, a ReLU activation function layer, a dropout layer, and a fully-connected layer. Then we feed the output to a softmax layer to predict the corresponding label. The training objective is to minimize the cross-entropy loss with  $\mathcal{L}_2$  regularization. Specifically, the optimal parameters  $\theta$  are obtained from

$$\mathcal{L}(\theta) = - \sum_{i=1}^n \hat{y}_i \log y_i + \lambda \sum_{\theta \in \Theta} \theta^2 \quad (3)$$

where  $\lambda$  is the regularization constant and  $\hat{y}_i$  is the predicted label corresponding to ground truth label  $y_i$ .

When no input transformation is used, the model is aspect-general and named as PLM-MeanPool and PLM-MeanPool-Concat for SC and OE, respectively. By incorporating the three input transformations, the model becomes more aspect-specific, denoted as +AC (Aspect Companion), +AP (Aspect Prompt), and +AM (Aspect Marker) respectively.

# 4 Experiments

## 4.1 Datasets

### 4.1.1 SC Datasets

Following previous work (Ma et al., 2021), we conduct experiments on two SC benchmarks to evaluate our models’ effectiveness and robustness. One is SemEval 2014 (Pontiki et al., 2014) (SEMEVAL), which contains data from laptop (SEM-LAP) and

restaurant (SEM-REST) domains; the other is the Aspect Robustness Test Set (ARTS-SC) (Xing et al., 2020), which is derived from the SEMEVAL dataset. Instances in ARTS-SC are generated with three adversarial strategies. These strategies enrich the test set from 638 to 1,877 for the laptop domain (ARTS-SC-LAP), and from 1,120 to 3,530 for the restaurant domain (ARTS-SC-REST). Note that each domain from SEMEVAL consists of separate training and test sets, while each domain from ARTS-SC only contains a test set. Since datasets in SEMEVAL do not contain development sets, 150 instances from the training set in each dataset are randomly selected to form the development set. Table 1 shows the statistics of the SC datasets.

#### 4.1.2 OE Datasets

For datasets used in OE (Fan et al., 2019; Wu et al., 2020), the original SEMEVAL benchmark annotates the aspects, but not the corresponding opinion spans, for each sentence. To solve the problem, (Fan et al., 2019) annotates the corresponding opinion spans for each given aspect in a sentence and removes the cases without explicit opinion spans. We use this variant in our OE experiments.

Since there is currently no robustness test set for OE, we follow (Xing et al., 2020)’s three adversarial strategies to generate an Aspect Robustness Test Set with spans (ARTS-OE) based on SEMEVAL. Specifically, we use these strategies to generate 1002 test instances for the laptop domain (ARTS-OE-LAP) and 2009 test instances for the restaurant domain (ARTS-OE-RES). Each aspect in a sentence is associated with an opinion span for OE. It is worth noting that this adversarial dataset can also be used for other tasks, e.g., aspect sentiment triplet extraction (Peng et al., 2020). Table 2 shows the statistics of the OE datasets. And the details of ARTS-OE are provided in Tabel 7. Since these OE datasets do not come with a development set, we randomly split 20% of the training set as validation set.

Dataset		#pos.	#neu.	#neg.
SEM-LAP	train	930	433	800
	test	341	169	128
	dev	57	27	66
SEM-REST	train	2,094	579	779
	test	728	196	196
	dev	70	54	26
ARTS-SC-LAP	test	883	407	587
ARTS-SC-REST	test	1,953	473	1,104

Table 1: Statistics of SC datasets.

Dataset		#sentences	#aspects
SEM-LAP	train	1,158	1,634
	test	343	482
SEM-REST	train	1,627	2,643
	test	500	865
ARTS-OE-LAP	test	1,002	2,404
ARTS-OE-REST	test	2,009	5,743

Table 2: Statistics of OE datasets. A sentence may contain multiple aspects. The number of aspect is identical to the number of pairs and instances.

## 4.2 Comparative Models and Baselines

We carry out an extensive evaluation of the proposed models (with and without input transformation), including *PLM-MeanPool* and *PLM-MeanPool +AC/AP/AM* for SC, *PLM-MeanPool-Concat* and *PLM-MeanPool-Concat +AC/AP/AM* for OE, against a wide range of baselines, categorized into two groups: PLM-based models and non-PLM models.

### 4.2.1 SC Baselines

*Non-PLM models* include: (a) IAN (Ma et al., 2017) interactively learns attentions between context words and aspect terms. (b) MemNet (Tang et al., 2016) applies attention multiple times on word memories, and the output of the last attention is used for prediction. (c) AOA (Huang et al., 2018) introduces an attention-over-attention based network to model interaction between aspects and contexts. (d) ASGCN (Zhang et al., 2019a) use graph convolutional networks to capture the aspect-specific information. *PLM-based models* include: (a) BERT/roBERTa-CLS-MLP use the representation of "[CLS]" as a classification feature to fine-tune the BERT/roBERTa model with an MLP layer. (b) AEN-BERT (Song et al., 2019) adopts BERT model and attention mechanism to model the relationship between contexts and aspects. (c) LCF-BERT (Zeng et al., 2019) employs Local-Context-Focus design with Semantic-Relative-Distance to discard unrelated sentiment words. (d) BERT/roBERTa-ASCNN is combined with BERT/roBERTa and ASCNN (Zhang et al., 2019a) model. (e) roBERTa-ASGCN (Zhang et al., 2019a) is combined with roBERTa and ASGCN.

### 4.2.2 OE Baselines

*Non-PLM models* include: (a) Pipeline (Fan et al., 2019) is a combination method of BiLSTM and Distance-rule method (Hu and Liu, 2004). (b) IOG (Fan et al., 2019) utilizes an Inward-Outward

LSTM and a Global LSTM to capture the information of aspect and global information, respectively. (c) LOTN Latent Opinions Transfer Network (Wu et al., 2020) uses an effective transfer learning method to identify latent opinions from the sentiment analysis model. (d) ARGCN (Jiang et al., 2021) is an extension of R-GCNs suited to encode syntactic dependency information to complete OE. *PLM-based models* include: (a) BERT+Distance-rule (Feng et al., 2021) is the combination of BERT and Distance-rule. (b) TF-BERT (Feng et al., 2021) utilizes the average pooling of target word embeddings to represent the target information, then it is fed into BERT to extract target-oriented opinion terms. (c) SDRN (Chen et al., 2020) utilizes BERT as the encoder, which consists of an opinion entity extraction unit, a relation detection unit, and a synchronization unit for the aspect opinion pair extraction task. (d) TSMSA-BERT (Feng et al., 2021) uses a target-specified sequence labeling method based on multi-head self-attention (TSMSA) to perform OE. (e) ARGCN+BERT (Jiang et al., 2021) adopts the last hidden states of the pretrained BERT as word representations and fine-tune it with the ARGCN model.

Implementation details and evaluation metrics can be found in Appendix A and B. It is worth noting that most previous methods did not use the dev set and may have overfitted the test set. We have made a systematic and comprehensive comparison for the first time under the same settings.

## 5 Results and Analysis

### 5.1 SC Results

Table 3 shows the standard (effectiveness) and robustness evaluation results for SC.

#### 5.1.1 Standard Results

Generally, our models with input transformations outperform the comparative baseline models. Before applying the transformations, our base models (BERT/RobERTa-MeanPool with aspect generality) perform equally good or even better than most baseline models.

Applying the input transformations, especially aspect marker (i.e., +AM), further improves performance significantly. For BERT-based models, the F1-scores of the BERT-MeanPool+AM model are 2.57% and 5.83% higher than AEN-BERT and LCF-BERT respectively on the SEM-REST standard dataset. For RobERTa-based models, the

three transformations are more effective. Specifically, the F1-scores of RobERTa-MeanPool+AC and RobERTa-MeanPool+AP improve by up to 1.54% and 1.28% on SEM-REST standard dataset. These results indicate that the proposed input transformations can promote PLMs to achieve effective aspect-specific context modeling.

Among the three transformations, in general AM performs better than AC and AP, indicating that AM is more effective for aspect-specific context modeling in PLMs. While the F1-scores of BERT-MeanPool+AM and RobERTa-MeanPool+AM gain improvements by 1.59% and 1.43% on SEM-REST, RobERTa-MeanPool+AM achieves the terrific results for SC, with F1-score are 78.5% and 79.58% on SEM-LAP and SEM-REST respectively.

#### 5.1.2 Robustness Results

We can see that the performances of the baseline models drop drastically on robustness test sets. In contrast, our models with the three transformations are more robust than the baseline models. The most effective and robust model is the RobERTa-MeanPool+AM, which achieves 72.59% and 74.04% of F1 score on the ARTS-SC-LAP and ARTS-SC-REST robustness test set, respectively, representing a 3.21% and 1.48% improvement over the strongest baseline RobERTa-ASGCN.

The three transformations significantly improve the BERT/RobERTa-MeanPool models' robustness, especially for RobERTa-MeanPool. Specifically, with AC, AP, and AM, the RobERTa-MeanPool model's F1-scores are improved by up to 1.30%, 1.36%, and 1.31% on ARTS-SC-REST robustness test set, respectively. The model with AM is more robust than the model with AC and AP. These robustness results demonstrate that the transformations can improve our models' robustness.

### 5.2 OE Results

Table 4 shows the standard and robustness results for OE.

#### 5.2.1 Standard Results

Before applying the transformations, our base models (BERT/RobERTa-MeanPool-Concat) perform poorly, even worse than most non-PLM baseline models. On the contrary, with the three transformations, our models perform significantly better than baseline models. Our BERT-based model with the three transformations achieves nearly

Models	SEM-LAP				SEM-REST			
	Standard		Robustness		Standard		Robustness	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
IAN	67.74	59.99	52.91	47.54	77.48	66.39	57.75	48.12
Memnet	67.81	60.67	52.00	46.50	76.77	64.46	55.30	46.67
AOA	69.47	63.13	52.00	46.50	77.57	66.02	58.19	49.02
ASGCN	70.97	65.31	56.59	52.12	78.87	68.12	64.89	55.41
AEN-BERT	77.37	71.83	71.49	66.37	83.66	75.50	73.24	66.31
LCF-BERT	76.55	71.40	71.19	66.95	81.66	72.24	70.57	62.75
BERT-CLS+MLP	75.42	69.08	54.91	51.21	78.95	67.66	53.86	47.16
RoBERTa-CLS+MLP	79.09	75.36	56.24	54.61	81.93	71.19	60.45	52.02
BERT-ASCNN	76.33	71.09	71.17	66.90	82.66	74.05	75.73	68.17
RoBERTa-ASCNN	81.41	77.22	73.59	70.14	85.93	78.01	78.85	70.69
RoBERTa-ASGCN	81.82	78.28	73.48	69.38	85.66	78.48	79.65	72.56
<b>BERT-MeanPool</b>	76.87	71.71	70.59	66.38	84.27	76.48	77.36	70.64
+AC	75.30	69.62	69.40	64.45	84.12	76.16	76.78	69.86
+AP	76.39	70.91	68.92	63.77	83.89	76.02	76.48	69.34
+AM	76.33	<b>71.93</b> ↑0.22	<b>70.78</b>	<b>67.06</b> ↑0.68	<b>84.71</b>	<b>78.07</b> ↑1.59	<b>78.10</b>	<b>72.38</b> ↑1.74
<b>RoBERTa-MeanPool</b>	81.38	77.68	74.67	71.21	85.41	78.15	79.75	72.73
+AC	<b>81.54</b>	77.54	<b>75.13</b>	71.02	<b>86.68</b> †	<b>79.69</b> †↑1.54	<b>80.63</b>	<b>74.03</b> †↑1.30
+AP	<b>81.85</b>	<b>77.91</b> ↑0.23	74.53	70.48	<b>86.43</b>	<b>79.43</b> †↑1.28	<b>80.72</b>	<b>74.09</b> †↑1.36
+AM	<b>82.07</b> †	<b>78.50</b> †↑0.82	<b>75.90</b> †	<b>72.59</b> †↑1.38	<b>86.41</b>	<b>79.58</b> †↑1.43	<b>80.88</b> †	<b>74.04</b> †↑1.31

Table 3: Standard and robust experimental results (%) on SC. The first and second blocks indicate non-PLM and PLM-based baseline models. Our models and better results are bold (Acc and F1, the larger, the better). The marker † represents that our models outperform the all other models significantly ( $p < 0.01$ ), and the small number next to each score indicates performance improvement (↑) compared with our aspect-general base model (BERT-MeanPool/RoBERTa-MeanPool).

identical results with the current state-of-the-art model (TSMSA-BERT). With AC, AP, and AM, the F1-scores of the RoBERTa-MeanPool-Concat model are improved by up to 13.04%, 12.89%, and 14.09% on SEM-LAP dataset, respectively. These results demonstrate that the three transformations can significantly promote PLMs to achieve effective aspect-specific context modeling for OE. Our RoBERTa-MeanPool-Concat+AM model achieves the new state-of-the-art result on OE.

### 5.2.2 Robustness Results

The performances of our base models (BERT/RoBERTa-MeanPool-Concat) drop drastically on robustness test set. Their F1-scores are only 39.68% and 38.76% on ARTS-OE-LAP and 44.23% and 56.93% on ARTS-OE-REST. In contrast, with the transformations, our models are more robust, achieving F1 scores up to 73.69% (RoBERTa-MeanPool-Concat+AM) on ARTS-OE-LAP, and 71.61% (RoBERTa-MeanPool-Concat+AP) on ARTS-OE-REST, demonstrating that the transformations can significantly improve our model’s robustness for OE.

## 5.3 Ablation Study

To further investigate the effects of the feature induction and the input transformations on aspect-specific context modeling of PLMs, we conduct ex-

tensive ablation experiments on standard datasets, whose results are included in Table 5 and 6 for SC and OE, respectively.

### 5.3.1 Aspect-specific Feature Induction

For SC and OE, we start with a simple base model that does not use the aspect feature induction component, but using just a context modeling representation after PLM and append an MLP layer. The base model is named as BERT/RoBERTa-CLS-MLP for SC, and BERT/RoBERTa-MLP for OE. Now we see what happens if we add back the aspect feature induction. For SC, our BERT/RoBERTa-MeanPool models always give a superior performance than the base model. The F1-scores of BERT-MeanPool are 2.63% and 8.82% higher than BERT-CLS-MLP on SEM-LAP and SEM-REST respectively. For OE, our BERT/RoBERTa-MeanPool-Concat models perform better than BERT/RoBERTa-MLP models. These results demonstrate the effectiveness of the aspect-specific feature induction methods with PLMs.

### 5.3.2 Aspect-specific Context Modeling

To investigate the effect of the aspect-specific context modeling with transformations, we add the input transformations to the above simple base models. The results show that the transformations bring significant performance improvements, even

Models	SEM-LAP		SEM-REST	
	Standard	Robustness	Standard	Robustness
Pipeline*	63.83	-	69.18	-
IOG*	70.99	-	80.23	-
LOTN*	72.02	-	82.21	-
ARGCN*	75.32	-	84.65	-
BERT+Distance-rule*	70.54	-	76.23	-
TF-BERT*	72.26	-	78.23	-
SDRN*	80.24	-	83.53	-
TMSA-BERT*	82.18	-	86.37	-
ARGCN-BERT*	76.36	-	85.42	-
<b>BERT-MeanPool-Concat</b>	68.27	39.68	69.08	44.23
+AC	80.31↑ <b>12.04</b>	70.98↑ <b>31.30</b>	85.09↑ <b>16.01</b>	70.01↑ <b>25.78</b>
+AP	79.60↑ <b>11.33</b>	68.06↑ <b>28.38</b>	85.32↑ <b>16.24</b>	70.25↑ <b>26.02</b>
+AM	81.06↑ <b>12.79</b>	71.23↑ <b>31.55</b>	85.62↑ <b>16.54</b>	69.68↑ <b>25.45</b>
<b>RoBERTa-MeanPool-Concat</b>	69.74	38.76	79.03	56.93
+AC	82.78↑ <b>13.04</b>	71.26↑ <b>32.50</b>	86.03↑ <b>7.00</b>	71.42↑ <b>14.49</b>
+AP	82.63↑ <b>12.89</b>	71.46↑ <b>32.30</b>	<b>86.58</b> ↑ <b>7.55</b>	<b>71.61</b> ↑ <b>14.68</b>
+AM	<b>83.83</b> ↑ <b>14.09</b>	<b>73.69</b> ↑ <b>34.93</b>	86.33↑ <b>7.30</b>	71.50↑ <b>14.57</b>

Table 4: Standard and robustness evaluation results (F1-score, %) on OE. The first and second blocks show the results of the non-PLM and BERT-based baseline models (with \*) respectively, which are extracted from the published papers (Wu et al., 2020) and (Feng et al., 2021). Note that there were no robustness results of the baseline models in the original published papers, so that we leave them blank. The results of our models are presented in the third and fourth blocks. The best results are bold (F1-score, the larger, the better).

Models	SEM-LAP	SEM-REST
BERT-MeanPool	71.71	76.48
BERT-CLS+MLP	69.08	67.66
+AC	68.82	74.03↑ <b>6.37</b>
+AP	70.47↑ <b>1.39</b>	76.78↑ <b>9.12</b>
+AM	70.24↑ <b>1.16</b>	74.19↑ <b>6.53</b>
RoBERTa-MeanPool	77.68	78.15
RoBERTa-CLS+MLP	75.36	71.19
+AC	77.62↑ <b>2.26</b>	76.04↑ <b>4.85</b>
+AP	78.40↑ <b>3.04</b>	78.53↑ <b>7.34</b>
+AM	78.21↑ <b>2.85</b>	79.91↑ <b>8.72</b>

Table 5: SC ablation experimental results (F1-score, %).

Models	SEM-LAP	SEM-REST
BERT-MeanPool-Concat	68.27	69.08
BERT-MLP	67.67	61.40
+AC	79.95↑ <b>12.28</b>	79.46↑ <b>18.06</b>
+AP	80.08↑ <b>12.41</b>	81.02↑ <b>19.62</b>
+AM	81.50↑ <b>13.83</b>	80.02↑ <b>18.62</b>
RoBERTa-MeanPool-Concat	69.74	79.03
RoBERTa-MLP	67.92	60.00
+AC	82.18↑ <b>14.26</b>	81.59↑ <b>21.59</b>
+AP	81.96↑ <b>14.04</b>	81.04↑ <b>21.04</b>
+AM	83.42↑ <b>15.50</b>	80.81↑ <b>20.81</b>

Table 6: OE ablation experimental results (F1-score).

better than the models with aspect feature induction. Especially the base models with the transformations for OE achieve nearly identical results to BERT/RoBERTa-MeanPool-Concat with transformations. These excellent results demonstrate the effectiveness of the proposed transformations for context modeling, which indirectly explains that context modeling is more critical than aspect fea-

ture induction for ABSA.

Model	Example
AG	[CLS] The food is ta #sty but the service is bad ! [SEP]
AC	[CLS] The food is ta #sty but the service is bad ! [SEP] food [SEP]
AP	[CLS] The food is ta #sty but the service is bad ! The target aspect is food [SEP]
AM	[CLS] The <asp> food </asp> is ta #sty but the service is bad ! [SEP]
AG	[CLS] The food is ta #sty but the service is bad ! [SEP]
AC	[CLS] The food is ta #sty but the service is bad ! [SEP] service [SEP]
AP	[CLS] The food is ta #sty but the service is bad ! The target aspect is service [SEP]
AM	[CLS] The food is ta #sty but the <asp> service </asp> is bad ! [SEP]

Figure 3: **Attention visualization.** Gradient saliency maps (Simonyan et al., 2014) for the embedding of each word in the transformations under BERT. Underlined words are aspects and corresponding opinion spans.

## 5.4 Visualization of Attention

To understand the effect of the three transformations, we visualize the attention scores separately offered by our OE model (BERT-MeanPool-Concat) with the transformations, as shown in Figure 3. The four attention vectors have encoded quite different concerns in the token sequence. We can observe that before applying the transformations, the model may attend to more irrelevant words. On the contrary, AC, AP, and AM can promote our model to attend to aspect-specific context and capture the correct opinion spans, thus achieving aspect-specific context modeling in PLM.

## 6 Conclusions

In this paper, we propose three aspect-specific input transformations and methods to leverage these

transformations to promote the PLM to pay more attention to the aspect-specific context in two aspect-based sentiment analysis (ABSA) tasks (SC and OE). We conduct experiments with standard benchmarks for SC and OE, along with adversarial ones for robustness tests. Our models with aspect-specific context modeling achieve the state-of-the-art performance for OE and outperform various strong models for SC. The extensive experimental results and further analysis indicated that aspect-specific context modeling can enhance the performance of ABSA.

## Acknowledgements

This research was supported in part by Natural Science Foundation of Beijing (grant number: 4222036) and Huawei Technologies (grant number: TC20201228005).

## References

- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Shaowei Chen, Jie Liu, Yu Wang, Wenzheng Zhang, and Ziming Chi. 2020. Synchronous double-channel recurrent network for aspect-opinion pair extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6515–6524.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Zhifang Fan, Zhen Wu, Xinyu Dai, Shujian Huang, and Jiajun Chen. 2019. Target-oriented opinion words extraction with target-fused neural sequence labeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2509–2518.
- Yuhao Feng, Yanghui Rao, Yuyao Tang, Ninghua Wang, and He Liu. 2021. Target-specified sequence labeling with multi-head self-attention for target-oriented opinion words extraction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1805–1815.
- Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177.
- Binxuan Huang and Kathleen M Carley. 2018. Parameterized convolutional neural networks for aspect level sentiment classification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1091–1096.
- Binxuan Huang, Yanglan Ou, and Kathleen M Carley. 2018. Aspect level sentiment classification with attention-over-attention neural networks. In *International Conference on Social Computing, Behavioral-Cultural Modeling and Prediction and Behavior Representation in Modeling and Simulation*, pages 197–206. Springer.
- Junfeng Jiang, An Wang, and Akiko Aizawa. 2021. Attention-based relational graph convolutional network for target-oriented opinion words extraction. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1986–1997.
- Xin Li, Lidong Bing, Wai Lam, and Bei Shi. 2018. Transformation networks for target-oriented sentiment classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 946–956.
- Xin Li, Lidong Bing, Wenxuan Zhang, and Wai Lam. 2019. Exploiting bert for end-to-end aspect-based sentiment analysis. *W-NUT 2019*, page 34.
- Yunlong Liang, Fandong Meng, Jinchao Zhang, Jinan Xu, Yufeng Chen, and Jie Zhou. 2019. A novel aspect-guided deep transition model for aspect based sentiment analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5569–5580.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. 2017. Interactive attention networks for aspect-level sentiment classification. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4068–4074.
- Fang Ma, Chen Zhang, and Dawei Song. 2021. [Exploiting position bias for robust aspect sentiment classification](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*,

- pages 1352–1358, Online. Association for Computational Linguistics.
- Haiyun Peng, Lu Xu, Lidong Bing, Fei Huang, Wei Lu, and Luo Si. 2020. Knowing what, how and why: A near complete solution for aspect-based sentiment analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8600–8607.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Maria Pontiki, Haris Papageorgiou, Dimitrios Galanis, Ion Androutsopoulos, John Pavlopoulos, and Suresh Manandhar. 2014. Semeval-2014 task 4: Aspect based sentiment analysis. *SemEval 2014*, page 27.
- Timo Schick and Hinrich Schütze. 2021. It’s not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2339–2352.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR*.
- Youwei Song, Jiahai Wang, Tao Jiang, Zhiyue Liu, and Yanghui Rao. 2019. Attentional encoder network for targeted sentiment classification. *arXiv preprint arXiv:1902.09314*.
- Youwei Song, Jiahai Wang, Zhiwei Liang, Zhiyue Liu, and Tao Jiang. 2020. Utilizing bert intermediate layers for aspect based sentiment analysis and natural language inference. *arXiv e-prints*, pages arXiv–2002.
- Duyu Tang, Bing Qin, and Ting Liu. 2016. Aspect level sentiment classification with deep memory network. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 214–224.
- Hao Tang, Donghong Ji, Chenliang Li, and Qiji Zhou. 2020. Dependency graph enhanced dual-transformer structure for aspect-based sentiment classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6578–6588.
- Kai Wang, Weizhou Shen, Yunyi Yang, Xiaojun Quan, and Rui Wang. 2020. Relational graph attention network for aspect-based sentiment analysis. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3229–3238.
- Shuai Wang, Sahisnu Mazumder, Bing Liu, Mianwei Zhou, and Yi Chang. 2018. Target-sensitive memory networks for aspect sentiment classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 957–967.
- Yequan Wang, Minlie Huang, Xiaoyan Zhu, and Li Zhao. 2016. Attention-based lstm for aspect-level sentiment classification. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 606–615.
- Zhen Wu, Fei Zhao, Xin-Yu Dai, Shujian Huang, and Jiajun Chen. 2020. Latent opinions transfer network for target-oriented opinion words extraction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 9298–9305.
- Xiaoyu Xing, Zhijing Jin, Di Jin, Bingning Wang, Qi Zhang, and Xuan-Jing Huang. 2020. Tasty burgers, soggy fries: Probing aspect robustness in aspect-based sentiment analysis. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3594–3605.
- Hu Xu, Bing Liu, Lei Shu, and S Yu Philip. 2019. Bert post-training for review reading comprehension and aspect-based sentiment analysis. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2324–2335.
- Rohan Kumar Yadav, Lei Jiao, Ole-Christoffer Granmo, and Morten Goodwin. 2021. Human-level interpretable learning for aspect-based sentiment analysis. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI-21)*. AAAI.
- 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*, pages 5408–5414.
- Biqing Zeng, Heng Yang, Ruyang Xu, Wu Zhou, and Xuli Han. 2019. Lcf: A local context focus mechanism for aspect-based sentiment classification. *Applied Sciences*, 9(16):3389.
- Chen Zhang, Qiuchi Li, and Dawei Song. 2019a. Aspect-based sentiment classification with aspect-specific graph convolutional networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4560–4570.
- Chen Zhang, Qiuchi Li, and Dawei Song. 2019b. Syntax-aware aspect-level sentiment classification with proximity-weighted convolution network. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1145–1148.

## A Implementation Details

For fair comparison, we re-produce all baselines based on their open-source codes under the same settings. For all the non-PLM models, 300-dimensional GloVe vectors (Pennington et al., 2014) are leveraged to initialize the input embeddings. All parameters of models are initialized with uniform distributions. The learning rate is  $10^{-3}$ . The coefficient of the L2 regularization is  $10^{-5}$ . In case a model has hidden states, the dimensionality of hidden states is set to 300. For experiments with BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) as the input embeddings, we adopt the BERT-base-uncased<sup>2</sup> model and the RoBERTa-base<sup>3</sup> model as our backbone network respectively, where the dimensionality of hidden states is 768 and the learning rate is set to  $10^{-5}$  for SC and  $5*10^{-5}$  for OE, while the regularization is removed. During all experiments, AdamW (Loshchilov and Hutter, 2019) is adopted optimizer in our models. The batch size is 64, and the maximal sequence length is 128. If a model involves attention mechanism, then the dot product-based attention is employed.

We also carry out experiments on two larger pre-trained models, i.e., BERT-Large and RoBERTa-Large. The experimental results show that the performances are similar to that of BERT-base and RoBERTa-base. Due to space limitation, we do not release the results on BERT-Large and RoBERTa-Large.

It is worth noting that most previous methods did not use the dev set and may have overfitted the test set. We have made a systematic and comprehensive comparison for the first time under the same settings.

## B Evaluation Metrics

For standard performance evaluation, each model is trained, validated and tested on the standard datasets for SC and OE. For SC, we use accuracy and macro-averaged F1-score as performance metrics. Following the previous work (Fan et al., 2019), we adopt F1-score only as the evaluation metric for OE. An opinion extraction is considered correct only when the opinion span predicted is the same as the ground truth.

To evaluate a model’s robustness on SC, the model is trained on the standard SEMEVAL datasets and tested on the corresponding ARTS-SC testsets. For a model’s robustness on OE, the model is trained on the standard SEMEVAL datasets and tested on the corresponding ARTS-OE testsets.

Finally, the experimental results are obtained by averaging five runs with random initialization. It is worth noting that our goal is to verify the effectiveness of the proposed method rather than achieving the sota on SC and OE. Such a simple method can achieve an effectiveness close to sota.

---

<sup>2</sup><https://huggingface.co/bert-base-uncased>.

<sup>3</sup><https://huggingface.co/roberta-base>.

Generation Strategy	Target Aspect: Opinion	Other Aspect:Opinion	Example
Source: The original sample from the test set	works : well positive	apple OS : happy	Works well , and I am extremely happy to be back to an apple OS .
RevTgt: Reverse the sentiment of the target aspect	works : badly negative	apple OS : happy	Works badly , but I am extremely happy to be back to an apple OS .
RevNon: Reverse the sentiment of the non-target aspects with originally the same sentiment as target	works : well positive	apple OS : unhappy	Works well , but I am extremely happy to be back to an apple OS .
AddDiff: Add aspects with the opposite sentiment from the target aspect	works : well positive	apple OS : happy games : issue video chat : iffy	Works well , and I am extremely happy to be back to an apple OS , but games being the main issue . And the video chat is the only thing that is iffy about it .

Table 7: The example of using three adversarial strategies to generate the Aspect Robustness Test Set with spans (**ARTS-OE**) based on SEMEVAL. Specifically, we use these strategies to generate 1002 test instances for the laptop domain (**ARTS-OE-LAP**) and 2009 test instances for the restaurant domain (**ARTS-OE-RES**). Each aspect in a sentence is associated with an opinion span for OE.