

# A Unified Model for Opinion Target Extraction and Target Sentiment Prediction\*

Xin Li<sup>†</sup> Lidong Bing<sup>‡</sup> Piji Li<sup>§</sup> and Wai Lam<sup>†</sup>

<sup>†</sup>Department of Systems Engineering and Engineering Management,  
The Chinese University of Hong Kong

<sup>‡</sup> R&D Center Singapore, Machine Intelligence Technology,  
Alibaba DAMO Academy

<sup>§</sup> Tencent AI Lab, Shenzhen, China

<sup>†</sup>{lixin, wlam}@se.cuhk.edu.hk, <sup>‡</sup>l.bing@alibaba-inc.com <sup>§</sup>pijili@tencent.com

## Abstract

Target-based sentiment analysis involves opinion target extraction and target sentiment classification. However, most of the existing works usually studied one of these two sub-tasks alone, which hinders their practical use. This paper aims to solve the complete task of target-based sentiment analysis in an end-to-end fashion, and presents a novel unified model which applies a unified tagging scheme. Our framework involves two stacked recurrent neural networks: The upper one predicts the unified tags to produce the final output results of the primary target-based sentiment analysis; The lower one performs an auxiliary target boundary prediction aiming at guiding the upper network to improve the performance of the primary task. To explore the inter-task dependency, we propose to explicitly model the constrained transitions from target boundaries to target sentiment polarities. We also propose to maintain the sentiment consistency within an opinion target via a gate mechanism which models the relation between the features for the current word and the previous word. We conduct extensive experiments on three benchmark datasets and our framework achieves consistently superior results.

## Introduction

Target-Based Sentiment Analysis (TBSA) aims to detect the opinion targets explicitly mentioned in sentences and predict the sentiment polarities over the opinion targets (Liu 2012; Pontiki 2014). For example, in the sentence “*USB3 Peripherals are noticeably less expensive than the ThunderBolt ones*”, the user mentions two opinion targets, namely, “*USB3 Peripherals*” and “*ThunderBolt ones*”, and expresses positive sentiment over the first, and negative sentiment over the second.

Traditionally, this task can be broken into two sub-tasks, namely, opinion target extraction and target sentiment classification. The goal of opinion target extraction is to detect the opinion target mentions in the text, and it has been extensively studied (Qiu et al. 2011; Liu, Xu, and Zhao 2013;

Liu, Xu, and Zhao 2014; Liu, Joty, and Meng 2015; Yin et al. 2016; Wang et al. 2016a; Wang et al. 2017; He et al. 2017; Li and Lam 2017; Li et al. 2018b). The second sub-task, i.e. target sentiment classification, performs as a multiplier for the usefulness of the extracted target mentions, as it can predict the sentiment polarity of the given opinion targets. This sub-task has also received a lot of attention in recent years (Dong et al. 2014; Tang, Qin, and Liu 2016; Wang et al. 2016b; Ma et al. 2017; Chen et al. 2017; Tay, Luu, and Hui 2017; Ma, Peng, and Cambria 2018; Hazarika et al. 2018; Li et al. 2018a; Wang et al. 2018; Xue and Li 2018; He et al. 2018; Li et al. 2019). However, most existing methods solving the second sub-task assume that the target mentions are given, which limits their practical use. To sum up, all the above works aim at solving only one of the sub-tasks. In order to apply these existing methods in practical settings, i.e. not only extracting the targets, but also predicting the target sentiment, one typical way is to pipeline the methods of the two sub-tasks together.

As observed in some other tasks (Jing et al. 2003; Ng and Low 2004; Finkel and Manning 2009; Miwa and Sasaki 2014), if two sub-tasks have strong couplings (e.g. NER and relation extraction), a more integrated model is usually more effective than a pipeline solution. For the TBSA task, previous researchers have attempted two approaches to a more integrated solution (Mitchell et al. 2013; Zhang, Zhang, and Vo 2015). One approach is to make the models of the two sub-tasks jointly trained, which utilizes a set of target boundary tags (e.g., B, I, E, S and O) and a set of sentiment tags (e.g. POS, NEG, NEU). The “joint” row of Table 1 gives an example of the tagging scheme in this approach. Another approach is to totally dismiss the boundary of the two sub-tasks, which utilizes a set of specially-designed tags (we name it “unified tagging scheme”), namely, B- $\{POS, NEG, NEU\}$ , I- $\{POS, NEG, NEU\}$ , E- $\{POS, NEG, NEU\}$ , S- $\{POS, NEG, NEU\}$ , denoting the beginning of, inside of, end of, and single-word opinion target with positive, negative or neutral sentiment respectively, and O denoting NULL sentiment. An example is given in the “unified” row in Table 1. Unfortunately, these initial attempts did not result in a more integrated model that can outperform the pipeline approaches.

Although the importance of solving the complete TBSA

\*The work described in this paper was mainly done when Lidong Bing was an employee and Xin Li was an intern at Tencent AI Lab. It is substantially supported by grants from the Research Grant Council of the Hong Kong Special Administrative Region, China (Project Codes: 14203414) and the Direct Grant of the Faculty of Engineering, CUHK (Project Code: 4055093). Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Input	The	AMD	Turin	Processor	seems	to	always	perform	much	better	than	Intel	.
Joint	O	B	I	E	O	O	O	O	O	O	O	S	O
	O	POS	POS	POS	O	O	O	O	O	O	O	NEG	O
Unified	O	B-POS	I-POS	E-POS	O	O	O	O	O	O	O	S-NEG	O

Table 1: Tagging schemes used in the integrated approaches. “Joint” and “Unified” refers to joint and unified approaches respectively.

task remains significant, existing studies are relatively less and their findings (Mitchell et al. 2013; Zhang, Zhang, and Vo 2015), to some extent, discouraged other researchers to do further explorations. However, we think that research efforts should be paid to explore a more integrated model for solving this task, because its two sub-tasks are highly coupled together and the potential of a more integrated model is promising.

In this paper, we investigate the complete task of TBSA and design a novel unified framework to handle it in an end-to-end fashion. The proposed framework involves two stacked Recurrent Neural Networks (RNN). The upper one produces the final tagging results of the TBSA task based on the unified tagging scheme. The lower one performs an auxiliary prediction of target boundaries with the aim for guiding and providing the information to the upper RNN. Such design is based on the observation that under the unified tagging scheme, the span information is exactly identical to that under the boundary tagging scheme. Refer to the example in Table 1, if a word is at the beginning of a target mention under the boundary scheme, i.e. having the tag B, it should also be at the beginning under the unified scheme, i.e. having the tag B-POS. In order to explore such inter-scheme tag dependency, we propose to guide the prediction of the upper RNN for the complete TBSA task with the boundary prediction from the auxiliary task, corresponding to the lower RNN. Specifically, we design a component to encode the dependencies into a transition matrix and use the matrix to map the probability distribution of the boundary prediction to the unified tag space of the TBSA task. Then, we determine the proportions of the obtained boundary-based probability scores in the tagging decision and consolidate them with the probability scores from the upper RNN for final predictions.

We also propose to maintain the consistency of the sentiment of individual words within the same target mention based on a simple gate mechanism. The gate mechanism is designed to explicitly consolidate the features of the current word and the previous word. Since both of the gate here and the transition matrix above need to take reliable boundary prediction for performing well, improving the reliability of such prediction in the lower RNN is supposed to be useful for the complete TBSA task. Therefore, we introduce another component to estimate the potential of a word to be a target word. Note that as defined by the task (Pontiki 2014; Pontiki 2015; Pontiki 2016), an opinion target should always co-occur with opinion words, thus, the words close to the opinion words are more likely to be target words and we obtain additional supervision signals for refining boundary

information based on this assumption.

In the experiments, our framework outperforms the state-of-the-art methods and the strongest sequence taggers on several benchmark datasets. We conducted detailed ablation studies to quantitatively demonstrate the effectiveness of the designed components. With some case analysis, we show how our framework can handle some difficult cases with the help of the designed components.

## Our Proposed Framework

### Task Definition

We formulate the complete Target-Based Sentiment Analysis (TBSA) task as a sequence labeling problem, and employ a unified tagging scheme:  $\mathcal{Y}^S = \{B-POS, I-POS, E-POS, S-POS, B-NEG, I-NEG, E-NEG, S-NEG, B-NEU, I-NEU, E-NEU, S-NEU, O\}$ . Except O, each tag contains two parts of tagging information: the boundary of aspect mention, and the aspect sentiment. For example, B-POS denotes the beginning of a positive aspect mention, and S-NEG denotes a single-word negative aspect. For a given input sequence  $X = \{x_1, \dots, x_T\}$  with length  $T$ , our goal is to predict a tag sequence  $Y^S = \{y_1^S, \dots, y_T^S\}$ , where  $y_i^S \in \mathcal{Y}^S$ .

### Model Description

**Overview** As shown in Figure 1, on the top of two stacked RNNs with LSTM cells, our framework designs three tailor-made components, depicted in detail with the callouts, to explore three important intuitions in the task of TBSA. Specifically, the upper  $LSTM^S$  is for the complete TBSA task and it predicts the unified tags as output, while the lower  $LSTM^T$  is for the auxiliary task and predicts the boundary tags of target mentions. The boundary prediction from  $LSTM^T$  is used to guide  $LSTM^S$  to make better predictions over the unified tags for the complete task.

The three key components are named **Boundary Guidance** (BG) component, **Sentiment Consistency** (SC) component and **Opinion-Enhanced** (OE) Target Word Detection component. The BG component takes the advantages of the boundary information provided by the auxiliary task to guide the  $LSTM^S$  for predicting the unified tags more accurately. The SC component is empowered with a gate mechanism to explicitly integrate the features of the previous word into the current prediction, aiming at maintaining the sentiment consistency within a multi-word opinion target. In order to provide boundary information of higher quality, the OE component, following the observation that “opinion targets and opinion words always co-occur”, performs another auxiliary

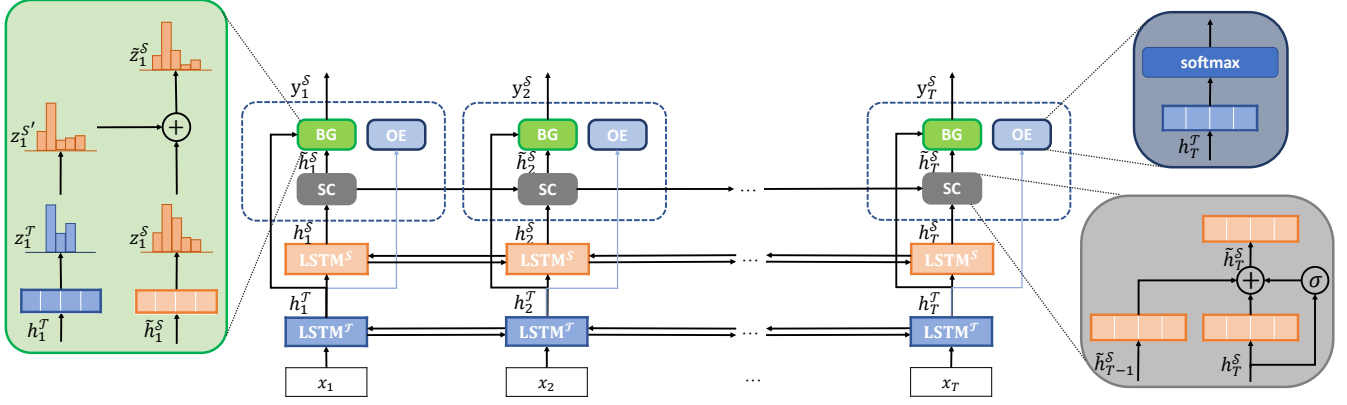


Figure 1: Architecture of the proposed framework.

binary classification task to determine if the current word is a target word.

**Target Boundary Guided TBSA** We employ  $LSTM^S$  with softmax decoding layer for the prediction of the tag sequence. It is observed that the boundary tag can provide important clues for the unified tag prediction. For example, if the current boundary tag is B, denoting the beginning of an opinion target, then the corresponding unified tag can only be B-POS, B-NEG or B-NEU. Thus, we introduce an additional network  $LSTM^T$  for the target boundary prediction, where the valid tag set  $\mathcal{Y}^T$  is  $\{B, I, E, S, O\}$ . We link these two LSTM layers so that the hidden representations generated by the  $LSTM^T$  can be directly fed to  $LSTM^S$  as guidance information. Specifically, their hidden representations  $h_t^T \in \mathbb{R}^{\dim_h^T}$  and  $h_t^S \in \mathbb{R}^{\dim_h^S}$  at the  $t$ -th time step ( $t \in [1, T]$ ) are calculated as follows:

$$\begin{aligned} h_t^T &= [\overrightarrow{LSTM^T}(x_t); \overleftarrow{LSTM^T}(x_t)], \\ h_t^S &= [\overrightarrow{LSTM^S}(h_t^T); \overleftarrow{LSTM^S}(h_t^T)], t \in [1, T]. \end{aligned} \quad (1)$$

The probability scores  $z_t^T \in \mathbb{R}^{|\mathcal{Y}^T|}$  over the boundary tags are calculated by a fully-connected softmax layer:

$$z_t^T = \mathbf{p}(y_t^T | x_t) = \text{Softmax}(\mathbf{W}^T h_t^T). \quad (2)$$

where the Softmax denotes the softmax activation function and  $\mathbf{W}^T$  is the model parameter. Similarly, the scores over the unified tags  $z_t^S \in \mathbb{R}^{|\mathcal{Y}^S|}$  are obtained as below:

$$z_t^S = \mathbf{p}(y_t^S | h_t^T) = \text{Softmax}(\mathbf{W}^S h_t^S). \quad (3)$$

As mentioned above, the boundary information is supposed to be useful for improving the performance of  $LSTM^S$ . (Zhang, Zhang, and Vo 2015) incorporated such boundary information by adding hard boundary constraints in the decoding step of the CRFs model. However, their prediction results are not promising. One reason is that their model employs a hard constraint which is prone to propagating the errors from the tagger of the boundary detection task and thus it decreases the performance of the TBSA tagger. Different from their way of imposing hard constraints,

our proposed BG component can absorb the boundary information via boundary guided transition and automatically determine its proportions in the final tagging decision based on the confidence of the target boundary tagger. Firstly, the BG component encodes the constraints into a transition matrix  $\mathbf{W}^{tr} \in \mathbb{R}^{|\mathcal{Y}^T| \times |\mathcal{Y}^S|}$ . As we have no prior knowledge about the transition probabilities between the boundary tags and the unified tags, we initially set them equally as follows:

$$\mathbf{W}_{i,j}^{tr} = \begin{cases} \frac{1}{|\mathcal{B}_i|}, & \text{if } j \in \mathcal{B}_i \\ 0, & \text{Otherwise} \end{cases} \quad (4)$$

where  $\mathcal{B}_i$  is the set of valid unified tags coherent with the boundary tag  $i$ . In this transition matrix, a non-zero element, e.g.,  $\mathbf{W}_{B, B-POS}^{tr}$ , denotes the probabilities of the unified tags given the boundary tag, and a zero element, e.g.,  $\mathbf{W}_{B, I-NEG}^{tr}$ , suggests that the unified tag cannot be inferred through this transition. After encoding the constraints, the next step is to guide the unified tag prediction with the boundary information. We directly propagate such information to the TBSA tagger by mapping the probability scores of the boundary tag  $z_t^T$  to the unified tag space. The transition-based sentiment score  $z_t^{S'} \in \mathbb{R}^{|\mathcal{Y}^S|}$  is obtained as follows:

$$z_t^{S'} = (\mathbf{W}^{tr})^\top z_t^T \quad (5)$$

where the transition operation is equivalent to the linear combination of the row vectors in the transition matrix  $\mathbf{W}^{tr}$ . Assuming  $z_t^T = [1, 0, 0, 0, 0]$  (i.e. taking the tag B), the result of the transition is exactly the row vector  $\mathbf{W}_{B, \cdot}^{tr}$ . As the unified tag can be partially derived from the boundary tag, a natural question is how to determine the proportions of the transition-based unified tagging scores  $z_t^{S'}$ . Intuitively, if the target boundary score  $z_t^T$  is nearly uniform, suggesting that the boundary tagger is not confident to its prediction, the obtained distribution over the unified tags, i.e.,  $z_t^{S'}$ , will also be close to a uniform distribution and has little meaningful information for the sentiment prediction. To avoid such uninformative boundary transitions, we calculate a proportion

score  $\alpha_t \in \mathbb{R}$  based on the confidence  $c_t$  of the target boundary tagger:

$$c_t = (z_t^T)^\top z_t^T$$

$$\alpha_t = \epsilon c_t \quad (6)$$

where the hyper-parameter  $\epsilon$  denotes the maximum proportions that the boundary-based scores  $z_t^{S'}$  occupy in the tagging decision. Obviously,  $c_t$  will be down-weighted if the boundary scores are uniformly distributed. The maximum confidence value is reached if  $z_t^T$  is a one-hot vector. The final scores are obtained by combining the boundary-based and model-based unified tagging scores:

$$\tilde{z}_t^S = \alpha_t z_t^{S'} + (1 - \alpha_t) z_t^S. \quad (7)$$

**Maintaining Sentiment Consistency** In the traditional target sentiment classification task, the sentiments towards the different words in a given multi-word opinion target are assumed to be identical. However, in the complete TBSA task, such sentiment consistency is not guaranteed since the task is formulated as a sequence tagging/labeling problem. Taking the sentence in Table 1 as an example, there is still some possibility that the word ‘‘Processor’’ is labeled with an E-NEG tag due to the independent tagging decisions made by LSTMs. To maintain the sentiment consistency within the same opinion target, we propose to predict the current unified tag using both of the features from the current and the previous time steps. Specifically, we design a Sentiment Consistency (SC) component with a gate mechanism to combine these two feature vectors:

$$\tilde{h}_t^S = g_t \odot h_t^S + (1 - g_t) \odot \tilde{h}_{t-1}^S$$

$$g_t = \sigma(\mathbf{W}^g h_t^S + \mathbf{b}^g) \quad (8)$$

where  $\mathbf{W}^g$  and  $\mathbf{b}^g$  are learnable parameters of the SC component, and  $\odot$  denotes the element-wise multiplication.  $\sigma$  is the sigmoid function. Through the gating, the previous features are considered in the current predictions and such indirect bi-gram dependency can help reduce the probability that the words within the same target hold different sentiments.

**Auxiliary Target Word Detection** A good boundary tagger for opinion targets is crucial for producing the boundary information of high quality. Here, we introduce the **OE** component to learn a more robust boundary tagger from another view of the training data. As defined in (Pontiki 2014; Pontiki 2015; Pontiki 2016), opinion targets are always collocated with opinion words. Inspired by this, we regard the word as a target word if there is at least one opinion word within the context window of fixed-size  $s$  of this word. Then, we train an auxiliary token-level classifier for discriminating target words and non-target words based on the distantly supervised labels and the boundary representations  $h_t^T$  are further refined with such supervision signals. The computational process of the **OE** component is below:

$$z_t^O = \text{Softmax}(\mathbf{W}^o h_t^T)$$

$$y_t^O = \arg \max_y z_t^O \quad (9)$$

where  $\mathbf{W}^o$  is the model parameter.

Dataset		Train	Dev	Test	Total
$\mathbb{D}_L$	# POS	883	104	339	1326
	# NEG	754	106	130	990
	# NEU	404	46	165	615
$\mathbb{D}_R$	# POS	2337	270	1524	4131
	# NEG	942	93	500	1535
	# NEU	614	50	263	927
$\mathbb{D}_T$	# POS	-			692
	# NEG	-			263
	# NEU	-			2244

Table 2: Statistics of the datasets.

## Model Training

All the components in our framework are differentiable, thus, the whole framework can be efficiently trained with gradient-based methods. Word/Token-level cross-entropy error is employed as the loss function:

$$\mathcal{L}^{\mathcal{I}} = -\frac{1}{T} \sum_{t=1}^T \mathbb{I}(y_t^{\mathcal{I},g}) \circ \log(z_t^{\mathcal{I}}) \quad (10)$$

where  $\mathcal{I}$  is the symbol of task indicator and its possible values are  $\mathcal{T}$ ,  $\mathcal{S}$ , and  $\mathcal{O}$ .  $\mathbb{I}(y)$  represents the one-hot vector with the  $y$ -th component being 1 and  $y_t^{\mathcal{I},g}$  is the gold standard tag for the task  $\mathcal{I}$  at the time step  $t$ . Then, the losses from the main TBSA task and the two auxiliary tasks are aggregated to form the training objective  $\mathcal{J}(\theta)$  of the framework:

$$\mathcal{J}(\theta) = \mathcal{L}^S + \mathcal{L}^T + \mathcal{L}^O. \quad (11)$$

## Experiments

### Dataset

Our model is evaluated on two product review datasets from SemEval ABSA challenges (Pontiki 2014; Pontiki 2015; Pontiki 2016) and the Twitter dataset. Table 2 gives the statistics of these benchmark datasets.  $\mathbb{D}_L$  (SemEval 2014) contains reviews from the laptop domain and the train-test split is the same as the original dataset.  $\mathbb{D}_R$  is the union set of the restaurant datasets from SemEval ABSA challenge 2014, 2015 and 2016. The new training dataset is obtained by merging the three years’ training datasets and the new testing set is built in the same way.  $\mathbb{D}_T$  consists of tweets collected by (Mitchell et al. 2013). The ground truth of the opinion target mentions and their sentiments are provided in these datasets. For  $\mathbb{D}_L$  and  $\mathbb{D}_R$ , we regard 10% randomly held-out training data as the development set. For  $\mathbb{D}_T$ , we report the ten-fold cross validation results, as done in (Mitchell et al. 2013; Zhang, Zhang, and Vo 2015), since there is no standard train-test split for this dataset.

The gold standard boundary annotations are available for the auxiliary target boundary prediction task. For another auxiliary task, namely, opinion-based target word detection, we employ the existing opinion lexicon<sup>1</sup> to provide the opinion words.

<sup>1</sup><http://mpqa.cs.pitt.edu/>

	Model	$\mathbb{D}_L$			$\mathbb{D}_R$			$\mathbb{D}_T$		
		P	R	F1	P	R	F1	P	R	F1
<b>Existing Baselines</b>	CRF-joint	57.38	35.76	44.06	60.00	48.57	53.68	43.09	24.67	31.35
	CRF-unified	59.27	41.86	49.06	63.39	57.74	60.43	48.35	19.64	27.86
	NN-CRF-joint	55.64	34.48	45.49	61.56	50.00	55.18	44.62	35.84	39.67
	NN-CRF-unified	58.72	45.96	51.56	62.61	60.53	61.56	46.32	32.84	38.36
<b>Pipeline Baselines</b>	CRF-pipeline	59.69	47.54	52.93	52.28	51.01	51.64	42.97	25.21	31.73
	NN-CRF-pipeline	57.72	49.32	53.19	60.09	61.93	61.00	43.71	37.12	40.06
	HAST-TNet	56.42	54.20	55.29	62.18	73.49	67.36	46.30	49.13	47.66
<b>Unified Baselines</b>	LSTM-unified	57.91	46.21	51.40	62.80	63.49	63.14	51.45	37.62	43.41
	LSTM-CRF-1	58.61	50.47	54.24	66.10	66.30	66.20	51.67	44.08	47.52
	LSTM-CRF-2	58.66	51.26	54.71	61.56	67.26	64.29	53.74	42.21	47.26
	LM-LSTM-CRF	53.31	59.4	56.19	68.46	64.43	66.38	43.52	52.01	47.35
<b>OURS</b>	Base model	60.00	46.85	52.61	61.48	66.16	63.73	53.02	41.47	46.50
	Base model + <b>BG</b>	58.58	50.63	54.31	67.51	66.42	66.96	52.26	43.84	47.66
	Base model + <b>BG</b> + <b>SC</b>	58.95	53.00	55.81	63.95	69.65	66.68	53.12	43.60	47.79
	Base model + <b>BG</b> + <b>OE</b>	63.43	49.53	55.62	62.85	66.77	65.22	53.10	43.50	47.78
	Full model	61.27	54.89	<b>57.90<sup>‡,‡</sup></b>	68.64	71.01	<b>69.80<sup>‡,‡</sup></b>	53.08	43.56	<b>48.01<sup>‡</sup></b>

Table 3: Main results of the complete TBSA task. “Base model” refers to the stacked LSTMs. The markers <sup>‡</sup> and <sup>‡</sup> refer to our full model significantly outperforms **HAST-TNet** and **LM-LSTM-CRF** respectively.

The evaluation metric measures the standard precision (P), recall (R) and F1 score based on the **exact** match, which means that an output segment is considered to be correct only if it exactly matches with the gold standard span of the target mention and the corresponding sentiment.

### Compared Models

We compare our framework with the following methods:

- **CRF-{pipeline, joint, unified}** (Mitchell et al. 2013): Conditional Random Fields (CRF) based sequence tagger<sup>2</sup>. “pipeline” denotes the pipeline approach. “joint” and “unified” are the models following the joint tagging scheme and unified tagging scheme respectively.
- **NN-CRF-{pipeline, joint, unified}** (Zhang, Zhang, and Vo 2015): Enhanced CRF models<sup>3</sup> armed with word embeddings and neural network feature extractors.
- **HAST-TNet**: HAST (Li et al. 2018b) and TNet (Li et al. 2018a) are the current state-of-the-art models on the tasks of target boundary detection and target sentiment classification respectively. HAST-TNet is the pipeline approach of these two models. We use the officially released codes<sup>4</sup> to produce the results.
- **LSTM-unified**: the standard LSTM model adopting the unified tagging scheme.
- **LSTM-CRF-1** (Lample et al. 2016): LSTM model with CRF decoding layer and no feature engineering is needed. We run the officially released code<sup>5</sup> and utilize the unified tag set to reproduce the results.

<sup>2</sup><http://www.m-mitchell.com/code/index.html>

<sup>3</sup><https://github.com/SUTDNLP/NNTargetedSentiment>

<sup>4</sup>Available at: <https://github.com/lixin4ever/HAST> and <https://github.com/lixin4ever/TNet> respectively.

<sup>5</sup><https://github.com/glample/tagger>

- **LSTM-CRF-2** (Ma and Hovy 2016): LSTM-CRF-2 is similar to LSTM-CRF-1. The difference is that LSTM-CRF-2 employs CNN rather than LSTM to learn the character-level word representations. We run the released code<sup>6</sup> to reproduce the results.
- **LM-LSTM-CRF** (Liu et al. 2018): Language model enhanced LSTM-CRF model. It is a competitive model in several sequence tagging tasks. We rerun their code<sup>7</sup> and report the tagging results based on the unified tagging scheme.

### Experiment Settings

**Word Embeddings** We use GloVe.840B.300d<sup>8</sup> released by (Pennington, Socher, and Manning 2014) to initialize the word embeddings, fine-tuned during training. The embeddings of the out-of-vocabulary words are sampled from the uniform distribution  $\mathcal{U}(-0.25, 0.25)$  (Kim 2014).

**Weight Initializations** The weight matrices in the LSTM units are initialized by following the Glorot Uniform strategy (Glorot and Bengio 2010) and the others are randomly sampled from the uniform distribution  $\mathcal{U}(-0.2, 0.2)$ . Besides, all biases are initialized as 0’s.

**Optimization** Our models are trained up to 50 epochs with Adam (Kingma and Ba 2014), with  $\beta_1 = \beta_2 = 0.9$ , and the initial learning rate  $\eta_0 = 10^{-3}$ . The decay rate is kept the same as the setting in (Lample et al. 2016). We apply dropout on word embeddings and the ultimate features for prediction. The dropout rates are empirically set as 0.5. The model obtaining the best F1 score on the development set is selected for producing the testing results.

<sup>6</sup><https://github.com/XuezheMax/NeuroNLP2>

<sup>7</sup><https://github.com/LiyuanLucasLiu/LM-LSTM-CRF>

<sup>8</sup><https://nlp.stanford.edu/projects/glove/>

**Others** Both of the dimension of the hidden representations  $\dim_h^T$  and  $\dim_h^S$  are 50. The maximum proportion  $\epsilon$  of the boundary-based scores is 0.5. The size of the context window  $s$  in the opinion-based target word detection component is 3. The tuning details of  $\epsilon$  and  $s$  are given later.

## Results and Analysis

**Main Results** Table 3 presents our comparisons with other methods for the complete TBSA task. To make the comparison fair, we use GloVe.840B.300d as the pre-trained word embeddings for all the baselines requiring word embedding input on all of the datasets. Besides, we align the train/dev/test configurations for all methods. The experimental results suggest that our proposed framework consistently gives the best F1 score across all datasets and significantly outperforms the strongest baselines in most cases.

Compared to HAST-TNet, the pipeline of two state-of-the-art models, our proposed framework achieves 2.6%, 2.4% and 0.40% absolute gains on  $\mathbb{D}_L$ ,  $\mathbb{D}_R$  and  $\mathbb{D}_T$  respectively, suggesting that a carefully-designed integrated model can be more effective than the pipeline approaches on the TBSA task. Three competitive unified sequence taggers (see the third block in Table 3) are also introduced into the comparative study. Again, our framework outperforms the best of them by 1.7%, 3.4% and 0.5% on the benchmark datasets. We notice that the improvement of our framework on the Twitter dataset is marginal in contrast with the unified baselines. The small gap is reasonable since these models employ additional component (e.g., LSTM or CNN) to learn the character-level word representations, whose capability for representing out-of-vocabulary words has been verified in (Santos and Zadrozny 2014; Kim et al. 2016), while our framework only utilizes the word-level features provided by the pre-trained word embeddings. Similar observation is captured in the comparison with HAST-TNet. We attribute this to the superior modeling power of the CNN applied in TNet when processing the ungrammatical sentences such as tweets and micro-blogs, as pointed out in (Li et al. 2018a).

We also notice that the performances of the CRF-based models, especially the recall (R) scores, are quite poor. Armed with the pre-trained word embeddings and neural network feature extractor, the models are slightly improved but the scores are still not promising.

**Effectiveness of the Proposed Components** To investigate the effectiveness of the designed components, we conduct ablation study on the proposed framework and the results are listed in the last block of the Table 3. Let us start the discussion from the base model, namely, the stacked LSTMs. We find that the base model always gives superior performance compared to the LSTM-unified. This result indicates that the boundary information predicted by the auxiliary LSTM indeed increases the F1 score of the complete TBSA task. With the help of the BG component, the performances are improved more significantly and the way we impose the boundary constraints proves effective for yielding more true positives. Another interesting finding is that

introducing the component **SC** or **OE** individually into the “Base model + **BG**” does not bring in too much gains on F1 measure and even hurts the prediction performance on  $\mathbb{D}_R$ . But putting them together, i.e., the “Full model”, leads to the new state-of-the-art result. This result illustrates the necessity of both of the **SC** and **OE** components in the boundary-guided TBSA. Considering the “Base model + **BG** + **SC**”, the quality of the boundary information may not be accurate without the clues from the **OE** component, and thus, the **SC** component tends to incorrectly align the sentiments of both the target words and non-target words. For the “Base model + **BG** + **OE**”, the quality of the boundary information obtained from the LSTM<sup>T</sup> is improved but the sentiments of the words within the same target are not fully consistent compared to the “Full model” armed with **SC** component. In summary, the **SC** component and the **OE** component are complementary to some extent when they are added into the boundary-guided “Base model + **BG**”.

**Case Analysis** Table 4 gives some prediction examples of the base model (i.e., the stacked LSTMs) and the models empowered with our proposed components. As observed in the first input and the second input, the “Base model” correctly predicts the target boundary but it fails to produce the right target sentiments, suggesting that linking the two LSTMs for the target boundary prediction and the TBSA task is still insufficient for exploiting the boundary information to improve the performance of the complete TBSA. The “Base model+**BG**” and the “Full model”, where the boundary constraints are properly imposed via our BG component, can correctly handle these two cases. Although the boundary information can guide the model to predict the sentiment more accurately, there is the possibility that only using the BG component (i.e. “Base model+**BG**”) inherits the errors from the lower boundary detection task, e.g., the third and the fourth input. Thus, the boundary information of high quality is crucial for improving the upper TBSA task and our OE component can serve as a simple but effective solution. Besides, we find that maintaining sentiment consistency within the same target mention, especially for those with several words (e.g., “*portobello and asparagus mole*” in the last input), is difficult for the “Base model” and “Base model+**BG**”, while our “Full model” alleviates this issue by employing the SC component to make predictions based on the features from the current and the previous time step.

**Impact of  $\epsilon$  and  $s$**  Here, we investigate the impacts of the maximum proportion  $\epsilon$  of the boundary-based scores and the window size  $s$  on the prediction performance. Specifically, the experiments are conducted on the development set of  $\mathbb{D}_R$ , the largest benchmark dataset. We vary  $\epsilon$  from 0.3 to 0.7, increased by 0.1, and two extreme values 0.0 and 1.0 are also included. The range of the window size  $s$  is 1 to 5. According to the results given in Figure 2, we observe that the best results are obtained at  $\epsilon=0.5$ . The  $\epsilon$  value basically affects the importance of the sentiment scores from the BG component in the final tagging decision and 0.5 is a good trade-off between absorbing boundary information and eliminating noises. We also observe that a moderate value of

Input	Base model		Base model + BG		Full model	
	Target	Complete	Target	Complete	Target	Complete
1. And the fact that it comes with an [ <i>i5 processor</i> ] <sub>POS</sub> definitely speeds things up	<i>i5 processor</i>	[ <i>processor</i> ] <sub>POS</sub> (✗)	<i>i5 processor</i>	[ <i>i5 processor</i> ] <sub>POS</sub>	<i>i5 processor</i>	[ <i>i5 processor</i> ] <sub>POS</sub>
2. There were small problems with [ <i>mac office</i> ] <sub>NEG</sub> .	<i>mac office</i>	[ <i>mac</i> ] <sub>NEG</sub> (✗)	<i>mac office</i>	[ <i>mac office</i> ] <sub>NEG</sub>	<i>mac office</i>	[ <i>mac office</i> ] <sub>NEG</sub>
3. The [ <i>teas</i> ] <sub>POS</sub> are great and all the [ <i>sweets</i> ] <sub>POS</sub> are homemade	<i>teas, sweets</i>	[ <i>teas</i> ] <sub>POS</sub> , [ <i>sweets</i> ] <sub>POS</sub>	<i>teas, sweets, homemade</i> (✗)	[ <i>teas</i> ] <sub>POS</sub> , [ <i>sweets</i> ] <sub>POS</sub> , [ <i>homemade</i> ] <sub>POS</sub> (✗)	<i>teas, sweets</i>	[ <i>teas</i> ] <sub>POS</sub> , [ <i>sweets</i> ] <sub>POS</sub>
4. I love the [ <i>form factor</i> ] <sub>POS</sub>	NONE	NONE	NONE	NONE	<i>form factor</i>	[ <i>form factor</i> ] <sub>POS</sub>
5. I blame the [ <i>Mac OS</i> ] <sub>NEG</sub> .	<i>Mac OS</i>	[ <i>Mac</i> <sub>NEG</sub> <i>OS</i> <sub>NEU</sub> ] <sub>NEG</sub> (✗)	<i>Mac OS</i>	[ <i>Mac</i> <sub>NEG</sub> <i>OS</i> <sub>POS</sub> ] <sub>NEG</sub> (✗)	<i>Mac OS</i>	[ <i>Mac OS</i> ] <sub>NEG</sub>
6. Also, I personally wasn't a fan of the [ <i>portobello and asparagus mole</i> ] <sub>NEG</sub> .	<i>portobello and asparagus mole</i>	[ <i>portobello</i> <sub>NEG</sub> and <sub>NEG</sub> <i>asparagus</i> <sub>NEG</sub> <i>mole</i> <sub>NEU</sub> ] <sub>NEG</sub> (✗)	<i>portobello and asparagus mole</i>	[ <i>portobello</i> <sub>NEG</sub> and <sub>NEG</sub> <i>asparagus</i> <sub>NEU</sub> <i>mole</i> <sub>NEU</sub> ] <sub>NEG</sub> (✗)	<i>portobello and asparagus mole</i>	[ <i>portobello and asparagus mole</i> ] <sub>NEG</sub>

Table 4: Case analysis. The “Target” column contains the results from the auxiliary task of target boundary detection. The “Complete” column presents the output of the complete TBSA task, but note that we only show the sentiment part of the unified labels (i.e. POS, NEG, and NEU) and use brackets to indicate the boundary. The marker ✗ denotes the incorrect prediction.

$s$  (i.e.,  $s = 3$ ) is the best for the TBSA task, probably because too large  $s$  may enforce the model to attend the larger context and increase the possibility of associating with irrelevant opinion words, on the other hand, too small  $s$  is likely not sufficient to involve the potential opinion words.

## Related Works

As mentioned in Introduction, Target-based Sentiment Analysis are usually divided into two sub-tasks, namely, the Opinion Target Extraction task (OTE) and the Target Sentiment Classification (TSC) task. Although these two sub-tasks are treated as separate tasks and solved individually in most cases, for more practical applications, they should be solved in one framework. Given an input sentence, the output of a method should contain not only the extracted opinion targets, but also the sentiment predictions towards them. Some previous works attempted to discover the relationship between these two sub-tasks and gave a more integrated solution for solving the complete TBSA task. Concretely, (Mitchell et al. 2013) employed Conditional Random Fields (CRF) together with hand-crafted linguistic features to detect the boundary of the target mention and predict the sentiment polarity. (Zhang, Zhang, and Vo 2015) further improved the performance of the CRF based method by introducing a fully connected layer to consolidate the linguistic features and word embeddings. However, they found that a pipeline method can beat both of the model with joint training and the unified model. In this paper, we reexamine the task, and proposed a new unified solution which outperforms all previous reported methods.

## Conclusions

We investigate the complete task of Target-Based Sentiment Analysis (TBSA), which is formulated as a sequence tagging problem with a unified tagging scheme in this paper. The basic architecture of our framework involves two stacked LSTMs for performing the auxiliary target boundary detection and the complete TBSA task respectively. On top of the base model, we designed two components

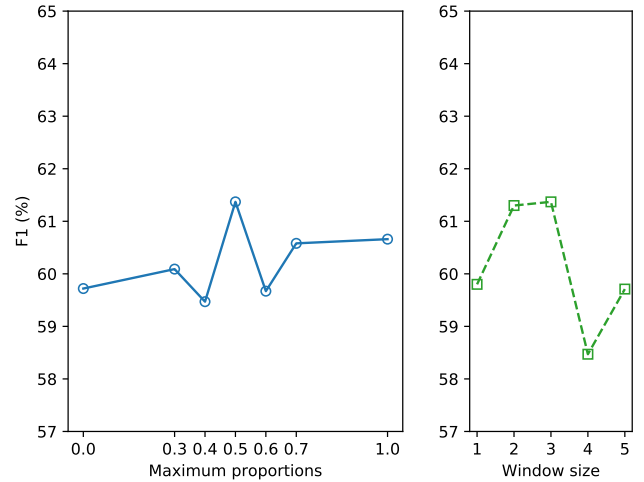


Figure 2: F1 scores (%) on the development set of  $\mathbb{D}_R$  with different  $\epsilon$  and  $s$  values.

to take the advantage of the target boundary information from the auxiliary task and maintain the sentiment consistency of the words within the same target. To ensure the quality of the boundary information, we employ an auxiliary opinion-based target word detection component to refine the predicted target boundaries. Experimental results and case studies well illustrate the effectiveness of our proposed framework, and a new state-of-the-art result of this task is achieved. We publicly release our implementation at <https://github.com/lixin4ever/E2E-TBSA>.

## References

- [Chen et al. 2017] Chen, P.; Sun, Z.; Bing, L.; and Yang, W. 2017. Recurrent attention network on memory for aspect sentiment analysis. In *EMNLP*, 452–461.
- [Dong et al. 2014] Dong, L.; Wei, F.; Tan, C.; Tang, D.; Zhou, M.; and Xu, K. 2014. Adaptive recursive neural net-



- work for target-dependent twitter sentiment classification. In *ACL*, 49–54.
- [Finkel and Manning 2009] Finkel, J. R., and Manning, C. D. 2009. Joint parsing and named entity recognition. In *NAACL-HLT*, 326–334.
- [Glorot and Bengio 2010] Glorot, X., and Bengio, Y. 2010. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 249–256.
- [Hazarika et al. 2018] Hazarika, D.; Poria, S.; Vij, P.; Krishnamurthy, G.; Cambria, E.; and Zimmermann, R. 2018. Modeling inter-aspect dependencies for aspect-based sentiment analysis. In *NAACL-HLT*, 266–270.
- [He et al. 2017] He, R.; Lee, W. S.; Ng, H. T.; and Dahlmeier, D. 2017. An unsupervised neural attention model for aspect extraction. In *ACL*, 388–397.
- [He et al. 2018] He, R.; Lee, W. S.; Ng, H. T.; and Dahlmeier, D. 2018. Exploiting document knowledge for aspect-level sentiment classification. In *ACL*, 579–585.
- [Jing et al. 2003] Jing, H.; Florian, R.; Luo, X.; Zhang, T.; and Ittycheriah, A. 2003. Howtogetachinesename(entity): Segmentation and combination issues. In *EMNLP*.
- [Kim et al. 2016] Kim, Y.; Jernite, Y.; Sontag, D.; and Rush, A. M. 2016. Character-aware neural language models. In *AAAI*.
- [Kim 2014] Kim, Y. 2014. Convolutional neural networks for sentence classification. In *EMNLP*, 1746–1751.
- [Kingma and Ba 2014] Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. In *ICLR*.
- [Lample et al. 2016] Lample, G.; Ballesteros, M.; Subramanian, S.; Kawakami, K.; and Dyer, C. 2016. Neural architectures for named entity recognition. In *NAACL-HLT*, 260–270.
- [Li and Lam 2017] Li, X., and Lam, W. 2017. Deep multi-task learning for aspect term extraction with memory interaction. In *EMNLP*.
- [Li et al. 2018a] Li, X.; Bing, L.; Lam, W.; and Shi, B. 2018a. Transformation networks for target-oriented sentiment classification. In *ACL*, 946–956.
- [Li et al. 2018b] Li, X.; Bing, L.; Li, P.; Lam, W.; and Yang, Z. 2018b. Aspect term extraction with history attention and selective transformation. In *IJCAI*.
- [Li et al. 2019] Li, Z.; Wei, Y.; Zhang, Y.; Xiang, Z.; and Li, X. 2019. Exploiting coarse-to-fine task transfer for aspect-level sentiment classification. In *AAAI*.
- [Liu et al. 2018] Liu, L.; Shang, J.; Xu, F.; Ren, X.; Gui, H.; Peng, J.; and Han, J. 2018. Empower sequence labeling with task-aware neural language model. In *AAAI*.
- [Liu, Joty, and Meng 2015] Liu, P.; Joty, S.; and Meng, H. 2015. Fine-grained opinion mining with recurrent neural networks and word embeddings. In *EMNLP*, 1433–1443.
- [Liu, Xu, and Zhao 2013] Liu, K.; Xu, L.; and Zhao, J. 2013. Syntactic patterns versus word alignment: Extracting opinion targets from online reviews. In *ACL*, 1754–1763.
- [Liu, Xu, and Zhao 2014] Liu, K.; Xu, L.; and Zhao, J. 2014. Extracting opinion targets and opinion words from online reviews with graph co-ranking. In *ACL*, 314–324.
- [Liu 2012] Liu, B. 2012. Sentiment analysis and opinion mining. *Synthesis lectures on human language technologies* 5(1):1–167.
- [Ma and Hovy 2016] Ma, X., and Hovy, E. 2016. End-to-end sequence labeling via bi-directional lstm-cnns-crf. In *ACL*, 1064–1074.
- [Ma et al. 2017] Ma, D.; Li, S.; Zhang, X.; and Wang, H. 2017. Interactive attention networks for aspect-level sentiment classification. In *IJCAI*, 4068–4074.
- [Ma, Peng, and Cambria 2018] Ma, Y.; Peng, H.; and Cambria, E. 2018. Targeted aspect-based sentiment analysis via embedding commonsense knowledge into an attentive lstm. In *AAAI*.
- [Mitchell et al. 2013] Mitchell, M.; Aguilar, J.; Wilson, T.; and Van Durme, B. 2013. Open domain targeted sentiment. In *EMNLP*.
- [Miwa and Sasaki 2014] Miwa, M., and Sasaki, Y. 2014. Modeling joint entity and relation extraction with table representation. In *EMNLP*.
- [Ng and Low 2004] Ng, H. T., and Low, J. K. 2004. Chinese part-of-speech tagging: One-at-a-time or all-at-once? word-based or character-based? In *EMNLP*.
- [Pennington, Socher, and Manning 2014] Pennington, J.; Socher, R.; and Manning, C. 2014. Glove: Global vectors for word representation. In *EMNLP*.
- [Pontiki 2014] Pontiki, M. 2014. Semeval-2014 task 4: Aspect based sentiment analysis. In *SemEval*, 27–35.
- [Pontiki 2015] Pontiki, M. e. a. 2015. Semeval-2015 task 12: Aspect based sentiment analysis. In *SemEval*, 486–495.
- [Pontiki 2016] Pontiki, M. e. a. 2016. Semeval-2016 task 5: Aspect based sentiment analysis. In *SemEval*, 19–30.
- [Qiu et al. 2011] Qiu, G.; Liu, B.; Bu, J.; and Chen, C. 2011. Opinion word expansion and target extraction through double propagation. *Computational Linguistics* 37(1).
- [Santos and Zadrozny 2014] Santos, C. D., and Zadrozny, B. 2014. Learning character-level representations for part-of-speech tagging. In *ICML*.
- [Tang, Qin, and Liu 2016] Tang, D.; Qin, B.; and Liu, T. 2016. Aspect level sentiment classification with deep memory network. In *EMNLP*.
- [Tay, Luu, and Hui 2017] Tay, Y.; Luu, A. T.; and Hui, S. C. 2017. Learning to attend via word-aspect associative fusion for aspect-based sentiment analysis. In *AAAI*.
- [Wang et al. 2016a] Wang, W.; Pan, S. J.; Dahlmeier, D.; and Xiao, X. 2016a. Recursive neural conditional random fields for aspect-based sentiment analysis. In *EMNLP*, 616–626.
- [Wang et al. 2016b] Wang, Y.; Huang, M.; zhu, x.; and Zhao, L. 2016b. Attention-based lstm for aspect-level sentiment classification. In *EMNLP*, 606–615.
- [Wang et al. 2017] Wang, W.; Pan, S. J.; Dahlmeier, D.; and



- Xiao, X. 2017. Coupled multi-layer attentions for co-extraction of aspect and opinion terms. In *AAAI*, 3316–3322.
- [Wang et al. 2018] Wang, S.; Mazumder, S.; Liu, B.; Zhou, M.; and Chang, Y. 2018. Target-sensitive memory networks for aspect sentiment classification. In *ACL*, 957–967.
- [Xue and Li 2018] Xue, W., and Li, T. 2018. Aspect based sentiment analysis with gated convolutional networks. In *ACL*, 2514–2523.
- [Yin et al. 2016] Yin, Y.; Wei, F.; Dong, L.; Xu, K.; Zhang, M.; and Zhou, M. 2016. Unsupervised word and dependency path embeddings for aspect term extraction. In *IJCAI*, 2979–2985.
- [Zhang, Zhang, and Vo 2015] Zhang, M.; Zhang, Y.; and Vo, D. T. 2015. Neural networks for open domain targeted sentiment. In *EMNLP*, 612–621.