

# Abstractive Opinion Tagging

Qintong Li\*  
Shandong University  
qintongli@mail.sdu.edu.cn

Piji Li  
Tencent AI Lab  
piji@tencent.com

Xinyi Li  
Peng Cheng Laboratory  
lixinyimichael@gmail.com

Zhaochun Ren†  
Shandong University  
zhaochun.ren@sdu.edu.cn

Zhumin Chen  
Shandong University  
chenzhumin@sdu.edu.cn

Maarten de Rijke  
University of Amsterdam  
& Ahold Delhaize  
m.derijke@uva.nl

## ABSTRACT

In e-commerce, *opinion tags* refer to a ranked list of tags provided by the e-commerce platform that reflect characteristics of reviews of an item. To assist consumers to quickly grasp a large number of reviews about an item, opinion tags are increasingly being applied by e-commerce platforms. Current mechanisms for generating opinion tags rely on either manual labelling or heuristic methods, which is time-consuming and ineffective. In this paper, we propose the *abstractive opinion tagging* task, where systems have to automatically generate a ranked list of opinion tags that are based on, but need not occur in, a given set of user-generated reviews.

The abstractive opinion tagging task comes with three main challenges: (1) the noisy nature of reviews; (2) the formal nature of opinion tags vs. the colloquial language usage in reviews; and (3) the need to distinguish between different items with very similar aspects. To address these challenges, we propose an abstractive opinion tagging framework, named AOT-Net, to generate a ranked list of opinion tags given a large number of reviews. First, a *sentence-level salience estimation* component estimates each review’s salience score. Next, a *review clustering and ranking* component ranks reviews in two steps: first, reviews are grouped into clusters and ranked by cluster size; then, reviews within each cluster are ranked by their distance to the cluster center. Finally, given the ranked reviews, a *rank-aware opinion tagging* component incorporates an alignment feature and alignment loss to generate a ranked list of opinion tags. To facilitate the study of this task, we create and release a large-scale dataset, called *eComTag*, crawled from real-world e-commerce websites. Extensive experiments conducted on the *eComTag* dataset verify the effectiveness of the proposed AOT-Net in terms of various evaluation metrics.

\*Work performed during an internship at Tencent AI Lab.

†Zhaochun Ren is the corresponding author.

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

WSDM '21, March 8–12, 2021, Virtual Event, Israel

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8297-7/21/03...\$15.00

<https://doi.org/10.1145/3437963.3441804>

## CCS CONCEPTS

• Information systems → Summarization; Sentiment analysis.

## KEYWORDS

Review analysis; abstractive summarization; e-commerce

### ACM Reference Format:

Qintong Li, Piji Li, Xinyi Li, Zhaochun Ren, Zhumin Chen, and Maarten de Rijke. 2021. Abstractive Opinion Tagging. In *Proceedings of the Fourteenth ACM International Conference on Web Search and Data Mining (WSDM '21)*, March 8–12, 2021, Virtual Event, Israel. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3437963.3441804>

## 1 INTRODUCTION

With the explosive growth of customer reviews in e-commerce scenarios, many online platforms, such as Amazon<sup>1</sup> and Alibaba<sup>2</sup>, provide opinion tags to enable potential buyers to make informed decisions without having to absorb large numbers of reviews. As shown in Figure 1, a sequence of opinion tags is a ranked list mined from a set of reviews that reflects different users’ preferences towards certain aspects of items. Many studies have focused on mining valuable information from reviews and shown promising results in various tasks, such as opinion summarization [1, 2, 5, 8, 24–26, 41] and item description generation [13, 36]. So far, however, no study seems to have developed opinion tagging methods to generate opinion tags that reflect diverse opinions of item aspects in a concise manner. In this paper, we propose the task of *abstractive opinion tagging*, which aims to automatically generate a ranked list of opinion tags that stem from, but need not occur in, a given set of user-generated reviews. This is an *abstractive* rather than an *extractive* opinion tagging task as the opinion tags do not need to occur in the review.

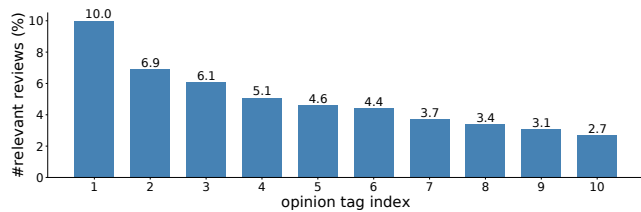
To solve this new task, we face three challenges. First, in reviews of e-commerce items, the noisy nature of the reviews inevitably makes it hard to identify salient item-related features [16, 24]. According to human annotations (on the eComTag dataset described below), we find that almost 59% of review sentences are not item-related or do not contain opinions towards certain aspects. Second, different reviewers have different ways of expressing themselves [44], whereas the target opinion tags are usually in a more formal style. The colloquial language usage of user reviews [49] makes abstractive methods necessary to learn better review representations. Third, it is difficult to reflect the different perspectives

<sup>1</sup><https://www.amazon.com/>

<sup>2</sup><https://www.alibaba.com/>

Reviews of a hot-pot restaurant	
$U_1$ :	The waitress was extremely attentive and even gave us a free fried man tou dessert that came with condensed milk for dipping...I love it!!
$U_2$ :	I was pleasantly surprised about how yummy the dish and the lamb were...
$U_3$ :	All in all; was a great experience and the service is really above and beyond.
$U_4$ :	The restaurant guest is more, can be served quickly, our table was quickly dish bowl filled with. Overall cool experience.
$U_5$ :	The shrimp was fresh and the pork mixture was tasty.
$U_6$ :	Fairly quick and polite service. It worth that price!
...	excellent, relaxed and cozy atmosphere, and what can I say, satisfying.
$U_N$ :	Food is delicious, reasonably priced...Go here! you deserve it!
Opinion tags	
hospitable service (223), delicious food (165), value for money (104), comfortable environment (65), served quickly (14).	

**Figure 1: An example of a set of reviews and their corresponding opinion tags, where  $U_N$  denotes the index of reviews, whereas the number behind each opinion tag reflects the number of reviews belonging to the tag.**

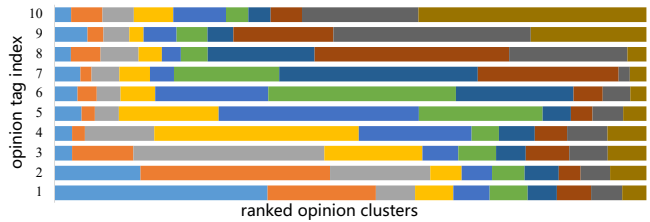


**Figure 2: Fraction of relevant reviews per opinion tag in the eComTag dataset.**

on an item in an accurate and diverse manner. Many e-commerce platforms *rank* opinion tags to help customers distinguish between different items with very similar aspects. In Figure 2, we plot the number of relevant reviews on different ranked indexes of opinion tags in the eComTag dataset (described below), which suggests a natural ranking for the tags when presented to users.

To address the challenges listed above, we design an abstractive framework, named AOT-Net, which consists of three components: (1) a *sentence-level salience estimation* component that predicts a salience score for each review; (2) a *review clustering and ranking* component that first groups reviews into clusters, which we refer to as “opinion clusters,” and ranks reviews by opinion cluster size; this component then ranks reviews within each cluster by their distance to the cluster center; and (3) a *rank-aware opinion tagging* component that generates opinion tags with ranks.

AOT-Net works in such a way that the ranks of the opinion clusters are correlated with the ranks of opinion tags. To see why this is potentially useful, we consider the eComTag dataset (crawled from e-commerce sites and consisting of items, reviews, and opinion tags), group reviews into opinion clusters and rank clusters by cluster size. We then compute the semantic similarity between the opinion clusters and the opinion tags, obtained by averaging the semantic similarity between an opinion tag and reviews in an opinion cluster. Figure 3 visualizes the average semantic similarity between 10 opinion tags and 10 opinion clusters for all item samples in the eComTag dataset. Clearly, the opinion tags have broader chunks with the opinion clusters at the same rank, revealing an alignment between the ranked opinion tags and ranked opinion clusters. Furthermore, the opinion tags are semantically similar



**Figure 3: The average semantic similarity between opinion tags and opinion clusters for all samples in the eComTag dataset. Each color denotes an opinion cluster. The same opinion clusters share the same color across different opinion tags. The width of the color band denotes the degree of semantic similarity between an opinion tag and the corresponding opinion cluster. (Best viewed in color.)**

to the neighbors of the corresponding opinion clusters, which is reflected by the relatively wide color bands, e.g., the second opinion tag is semantically similar to the first opinion cluster and the third opinion cluster. The *rank-aware opinion tagging* component of AOT-Net integrates two alignment strategies: (1) an opinion tag and its corresponding opinion clusters are placed at similar ranks; and (2) an *alignment loss* explicitly encourages the model to focus on the aligned opinion clusters and ignore others.

To validate AOT-Net, we collect a new dataset, named eComTag, from Chinese e-commerce websites, containing reviews and opinion tags for 50,068 items. Our experiments based on the eComTag dataset show that AOT-Net is capable of significantly improving the generation performance on abstractive opinion tagging over state-of-the-art baselines.

Our contributions can be summarized as follows:

- We propose a new task, *abstractive opinion tagging*, to generate opinion tags based on large volumes of item reviews.
- We propose an abstractive framework AOT-Net to generate opinion tags based on reviews for a given item. AOT-Net has a sentence-level salience estimation component and a review clustering and ranking component to highlight salient reviews and rank reviews by clustering. We propose a rank-aware opinion tagging component with two alignment strategies to generate ranked opinion tags.
- We collect a large-scale dataset, namely eComTag, consisting of item reviews and opinion tags, to support research into abstractive opinion tagging. Experimental results conducted on the eComTag dataset demonstrate the effectiveness of the AOT-Net framework.

## 2 RELATED WORK

Related work comes in two categories: keyphrase generation and opinion summarization.

### 2.1 Keyphrase Generation

A lot of research has been conducted on generating keyphrases to summarize various types of text such as tweets, news reports, research articles, etc. [11, 28, 29, 34, 37, 49, 50, 53]. Early approaches to keyphrase generation extract important phrases from the document

as the results. Sequence tagging models have been applied to identify keyphrases [18, 31, 52]. Retrieval-based approaches utilize a two-step pipeline to extract and rank candidate keyphrases [21, 33, 35, 48]. Sun et al. [42] adopt an extractive graph-based approach, which applies a point network to generate a set of diverse keyphrases. Recently, abstractive approaches have also been explored. [?] are the first to employ attention-based seq2seq framework with copy mechanism to conduct abstractive keyphrase generation. Chan et al. [10] propose a reinforcement learning approach for neural keyphrase generation that encourages a model to generate both sufficient and accurate keyphrases. Wang et al. [49] propose a topic-aware neural keyphrase generation method to identify topic words.

Unlike the work listed above, which only considers keyphrase generation for single document, we consider opinion tagging from multiple documents, that is, from all of the reviews for a given item.

## 2.2 Opinion Summarization

Opinion summarization has become an emerging research topic in recent years. Early studies on opinion summarization focus on extracting salient sentences from the original review text [3, 7, 15, 19, 30, 51]: Hu and Liu [19] identify item features mentioned in the reviews and then extract opinion sentences for the identified features. Xiong and Litman [51] utilize unsupervised learning methods to extract review summaries by exploiting review helpfulness ratings. Angelidis and Lapata [3] present a weakly supervised neural framework for aspect-based opinion summarization by combining the tasks of aspect extracting and sentiment predicting. Reflecting the most representative opinions from reviewers, many recent studies have shown that abstractive approaches are more appropriate for summarizing review text [6, 14, 17, 23, 45]: Gerani et al. [17] utilize a template filling strategy to indirectly generate a review summary; Wang and Ling [47] apply an attention-based encoder-decoder framework to generate an abstractive summary for opinionated documents. The main objective of the above summarization approaches is to generate coherent sentences to summarize opinions.

In contrast, we propose the *abstractive opinion tagging* task so as to generate opinion tags from a large number of user-generated reviews. In our scenario, opinion tags are more concise but without loss of essential information; they should help users comprehend reviews quickly and conveniently [26].

## 3 PROBLEM FORMULATION

Before detailing our proposed method, AOT-Net, we first formulate the abstractive opinion tagging problem. We use bold lowercase characters to denote vectors, and bold upper case characters to denote matrices. We write  $\mathbf{W}$  and  $\mathbf{b}$  for a projection matrix and a bias vector in a neural network layer, respectively. Suppose that there are  $M$  reviews for a given item. We denote each review  $X_i$ ,  $1 \leq i \leq M$  as a sequence of words, i.e.,  $X_i = [x_1, \dots, x_{L_{x_i}}]$ , where  $L_{x_i}$  denotes the number of words in  $X_i$ . In the same way, we assume that  $N$  opinion tags exist for a given item. We denote each opinion tag  $Y_j$ ,  $1 \leq j \leq N$  as a sequence of words, i.e.,  $Y_j = [y_1, \dots, y_{L_{y_j}}]$ , where  $L_{y_j}$  refers to the number of words in  $Y_j$ . Given a set of reviews  $\mathcal{X} = \{X_1, X_2, \dots, X_M\}$ , the task of *abstractive opinion tagging* is to generate a sequence of opinion tags  $\mathcal{Y} = [Y_1, Y_2, \dots, Y_N]$ .

## 4 METHOD

### 4.1 Overview

Before providing the details of AOT-Net, our proposed method for abstractive opinion tagging, we first provide an overview in Figure 4. We divide AOT-Net into three main phases: (A) sentence-level salience estimation; (B) review clustering and ranking; and (C) rank-aware opinion tagging. For a set of reviews  $\mathcal{X}$  about a given item, in phase A, we derive a salience score  $z_i$  for each review  $X_i \in \mathcal{X}$  to estimate its item-aware salience information. In phase B, reviews are first encoded into vector representations and weighted by corresponding salience scores. The weighted vector representations of reviews are clustered into  $K$  opinion clusters  $\{C_1, \dots, C_K\}$  and ranked by cluster size. Reviews within each cluster are ranked by their distance to the cluster center. Then we flatten ranked reviews into word-level vector representations. In phase C, we use review representations to generate ranked opinion tags via two alignment constraints, i.e., alignment features and alignment loss. We jointly learn all components in a multi-task learning framework.

### 4.2 A: Sentence-level Salience Estimation

The aim of the sentence-level salience estimation component is to compute a salience score for each review  $X_i \in \mathcal{X}$ . We design a *sentence-level self-attention mechanism* to highlight item-related reviews and reduce noise. First, the component reads each review sequence  $X_i = [x_1, \dots, x_{L_{x_i}}]$  and uses a lookup table to convert each review word  $x_p$  to a word embedding vector  $\mathbf{x}_p \in \mathbb{R}^{d_e}$ . To incorporate the contextual information of the review text into the representation of each word, we feed each embedding vector  $\mathbf{x}_p$  to a bi-directional Gated-Recurrent Unit (GRU) [12] to learn a hidden representation  $\mathbf{h}_p \in \mathbb{R}^d$ . More specifically, a bi-directional GRU consists of a forward GRU that reads the embedding sequence from  $x_1$  to  $x_{L_{x_i}}$  and a backward GRU that reads from  $x_{L_{x_i}}$  to  $x_1$ :

$$\vec{\mathbf{h}}_p = \text{GRU}_f(\mathbf{x}_p, \vec{\mathbf{h}}_{p-1}), \quad (1)$$

$$\overleftarrow{\mathbf{h}}_p = \text{GRU}_b(\mathbf{x}_p, \overleftarrow{\mathbf{h}}_{p+1}), \quad (2)$$

where  $\vec{\mathbf{h}}_p \in \mathbb{R}^{d/2}$  and  $\overleftarrow{\mathbf{h}}_p \in \mathbb{R}^{d/2}$  denote the hidden states of the forward GRU<sub>f</sub> and backward GRU<sub>b</sub>, respectively. We concatenate the last forward hidden state  $\vec{\mathbf{h}}_{L_{x_i}}$  and last backward hidden state  $\overleftarrow{\mathbf{h}}_1$  to form the hidden representation for review  $X_i$ , i.e.,  $\mathbf{h}_{X_i} = [\vec{\mathbf{h}}_{L_{x_i}}; \overleftarrow{\mathbf{h}}_1]$ .

Next, we pass hidden representations of all reviews  $\{\mathbf{h}_{X_i}\}_{i=1:M}$  to a self-attention layer to model more complex interactions among the reviews. We propose a salience context vector  $\mathbf{c}_i$  for each review  $X_i$  to denote the shared information from other reviews:

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{h}_{X_i}, \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{h}_{X_i}, \quad \mathbf{v}_i = \mathbf{W}_v \mathbf{h}_{X_i}, \quad (3)$$

$$\mathbf{c}_i = \sum_{i'=1}^M \frac{\exp(\mathbf{q}_i^T \mathbf{k}_{i'})}{\sum_{o=1}^M \exp(\mathbf{q}_i^T \mathbf{k}_o)} \mathbf{v}_{i'}, \quad (4)$$

where  $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i \in \mathbb{R}^{d \times d}$  refer to query, key, and value vectors, respectively. These vectors are linearly transformed from review hidden representation  $\mathbf{h}_{X_i}$ . Then we apply a residual connection from the review hidden representation  $\mathbf{h}_{X_i}$  to the salience context

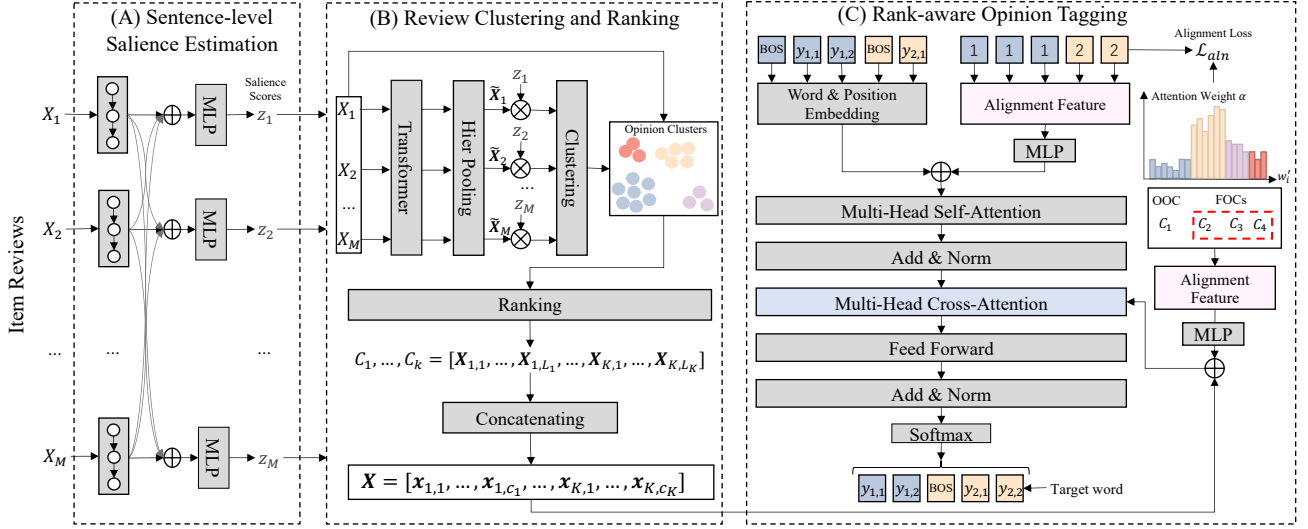


Figure 4: Our proposed framework AOT-Net for abstractive opinion tagging.

vector  $\mathbf{c}_i$  and feed it to a two-layer feed-forward network with a ReLU as the activation function:

$$\mathbf{h}'_{X_i} = \mathbf{W}_{s1}(\text{ReLU}(\mathbf{W}_{s2}(\mathbf{h}_{X_i} + \mathbf{c}_i))). \quad (5)$$

Given the context-enhanced review hidden representation  $\mathbf{h}'_{X_i}$ , we can derive the real-valued saliency score  $z_i$  for  $X_i$ :

$$z_i = \sigma(\mathbf{W}_s \mathbf{h}'_{X_i} + b_s), \quad (6)$$

where  $\mathbf{W}_s \in \mathbb{R}^d$  and  $b_s \in \mathbb{R}$ .  $\sigma(\cdot)$  is the sigmoid activation function. The saliency scores  $\{z_1, \dots, z_M\}$  serve as the saliency weights of review representations for the later review clustering and ranking component.

In order to optimize the sentence-level saliency estimating component, we manually label a binary *saliency label*  $z_i^* \in \{0, 1\}$  for each review  $X_i$ , where 1 denotes “item-related” whereas 0 denotes “noisy”. Then, sentence-level saliency estimation component is trained by minimizing the cross-entropy loss function:

$$\mathcal{L}_{cla} = -\frac{1}{M} \sum_{i=1}^M z_i^* \log(z_i) + (1 - z_i^*) \log(1 - z_i). \quad (7)$$

### 4.3 B: Review Clustering and Ranking

We propose a review clustering and ranking component to learn the ranks of reviews by grouping reviews into ranked opinion clusters, which is the main prerequisite to accurately generate ranked opinion tags.

We use a standard transformer encoder [46] to convert each review  $X_i$  into vector representations. Following Vaswani et al. [46], we first map each word  $x \in X_i$  into its vectorized representation  $\mathbf{x} \in \mathbb{R}^{d_e}$  using a word embedding layer and a positional embedding layer, as shown in the following equation:

$$\mathbf{x} = \text{Embed}(x) + \text{Pos}(x). \quad (8)$$

Then we use a transformer layer to encode global contextual information for words within  $X_i$ .

$$\mathbf{g} = \text{LayerNorm}(\mathbf{x}^{n-1} + \text{MHAtt}(\mathbf{x}^{n-1})), \quad (9)$$

$$\mathbf{x}^n = \text{LayerNorm}(\mathbf{g} + \text{FFN}(\mathbf{g})), \quad (10)$$

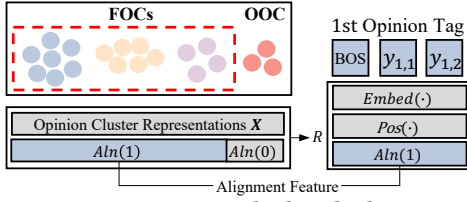
where LayerNorm is the layer normalization proposed by Ba et al. [4]; MHAtt is the multi-head attention mechanism introduced by Vaswani et al. [46]; FFN is a two-layer feed-forward network with ReLU as hidden activation function; and  $n$  is the number of transformer block layers. The word-level vector representations of review  $X_i$  are  $\mathbf{X}_i = [\mathbf{x}_1, \dots, \mathbf{x}_{L_{X_i}}] = [\mathbf{x}_1^n, \dots, \mathbf{x}_{L_{X_i}}^n]$ . To obtain the sentence representation  $\tilde{\mathbf{X}}_i$  of review  $X_i$ , we perform a hierarchical pooling operation [39] across its different words. The hierarchical pooling mechanism can preserve word order information and has demonstrated superior performance over mean-pooling or max-pooling on many semantic analysis tasks [39].

To highlight the item-aware reviews and ignore noisy reviews, the sentence representations of reviews are first weighted by the corresponding saliency scores, i.e.,  $\tilde{\mathbf{X}}'_i = z_i \tilde{\mathbf{X}}_i$ . Then we apply the  $k$ -means [32] algorithm on  $\{\tilde{\mathbf{X}}'_1, \dots, \tilde{\mathbf{X}}'_M\}$  to group corresponding  $\{\mathbf{X}_1, \dots, \mathbf{X}_M\}$  into  $K$  opinion clusters.<sup>3</sup> We rank opinion clusters from the largest (representing the highest number of reviews) to the smallest, denoted as  $[C_1, \dots, C_K]$ . For each cluster, we rank reviews from the nearest (representing the distance between review and cluster center) to the farthest. Finally, we obtain a ranked list of reviews, represented as  $[\mathbf{X}_{1,1}, \dots, \mathbf{X}_{1,L_1}, \dots, \mathbf{X}_{K,1}, \dots, \mathbf{X}_{K,L_K}]$ , where  $\mathbf{X}_{k,i}$  is the  $i$ -th review in the  $k$ -th opinion cluster and  $L_k$  is the number of reviews in the  $k$ -th opinion cluster.

We sequentially concatenate the vector representations of reviews in the ranked list and derive the final word-level representations, i.e.,  $\mathbf{X} = [\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,c_1}, \dots, \mathbf{x}_{K,1}, \dots, \mathbf{x}_{K,c_K}]$ .<sup>4</sup> Similarly,  $x_{k,p}$  is the  $p$ -th word in the  $k$ -th opinion cluster and  $c_k$  is the number of words in the  $k$ -th opinion cluster. Next,  $\mathbf{X}$  will serve as the memory bank for the later rank-aware opinion tagging component.

<sup>3</sup> $K$  is manually assigned according to the number of reviews. If  $M \leq 200$ ,  $K = \lceil \frac{M}{20} \rceil$ , otherwise,  $K = 20$ .

<sup>4</sup>In this paper, we only focus on the ranks of opinion clusters. We regard reviews or words in the same opinion cluster as equally important.



**Figure 5: Representations attached with alignment features for opinion tags and opinion clusters.**

#### 4.4 C: Rank-aware Opinion Tagging

Accurately generating opinion tags with ranks is challenging. Therefore, we propose a rank-aware opinion tagging component to generate ranked opinion tags.

In the training stage, we add a start token BOS at the beginning of each opinion tag, i.e.,  $Y'_j = [\text{BOS}; Y_j]$  where  $[\cdot]$  is the concatenation function. Then we concatenate all opinion tags into a sequence of words:  $Y = [Y'_1; \dots; Y'_N] = [y_{1,0}, \dots, y_{1,L_{y_1}}, \dots, y_{N,0}, \dots, y_{N,L_{y_N}}]$ , where  $y_{j,q}$  is the  $q$ -th word of opinion tag  $Y'_j$ ,  $y_{j,0}$  is the BOS token of the  $j$ -th opinion tag. Our decoder, i.e., the rank-aware opinion tagging component, follows the transformer architecture [46].

From the data analysis in Figure 3, we know that the ranks of opinion tags have a strong correlation with the ranks of opinion clusters. The  $j$ -th opinion tag may pay attention to the  $j$ -th opinion cluster and its surrounding neighbors simultaneously. Therefore, for each opinion tag  $Y'_j$ , we hypothesize that the model needs to focus on the  $F$  most related opinion clusters.<sup>5</sup> If an opinion cluster belongs to the  $F$  focused opinion clusters, we call it a *Focused Opinion Cluster* (FOC). Otherwise, we call it an *Outer Opinion Cluster* (OOC). We design two alignment strategies between FOCs and opinion tags, i.e., *incorporating alignment features* and *enforcing alignment loss*, which help to improve the generation of ranked opinion tags.

*Alignment Feature.* Intuitively, the opinion tags and their FOCs are semantically similar in the vector space. To help the model capture the alignment between opinion tags and their FOCs, we incorporate alignment features  $\text{Aln}(\cdot)$  into the word-level representations of the opinion tags and opinion clusters. Formally,  $\text{Aln}(\cdot)$  is a function that maps an integer into a vector. Now, we explain how we will use it to represent the ranks of opinion tags and opinion clusters. First, we incorporate the alignment feature into the representations of opinion tags. The words in the  $j$ -th opinion tag have the same rank  $j$  where  $1 \leq j \leq N$ . For each word  $y_{j,q}$ , the vectorized representation is the sum of the alignment feature  $\text{Aln}(j)$ , the word embedding  $\text{Embed}(y_{j,q})$ , and the positional embedding  $\text{Pos}(y_{j,q})$ :

$$y_{j,q} = \mathbf{W}_{rt} \text{Aln}(j) + \text{Embed}(y_{j,q}) + \text{Pos}(y_{j,q}), \quad (11)$$

where  $\mathbf{W}_{rt} \in \mathbb{R}^{d_e \times d_e}$  is a trainable model parameter.

For the target  $j$ -th opinion tag, the ranks of FOCs are set to  $j$  as well, while the ranks of OOCs are set to 0. Then we enhance the word vectors in  $\mathbf{X}$  with alignment features to capture the alignment

<sup>5</sup> $F$  is a hyperparameter.  $F$  opinion clusters mean the  $j$ -th opinion clusters and its surrounding neighbors.

between reviews and opinion tags:

$$\text{aln}_{x_{i,p}} = \begin{cases} \text{Aln}(j), & x_{i,p} \in \text{FOCs} \\ \text{Aln}(0), & x_{i,p} \in \text{OOCs}. \end{cases} \quad (12)$$

Next, the alignment features  $[\text{aln}_{x_{1,1}}, \dots, \text{aln}_{x_{K,c_K}}]$  are added into the review representations  $[\mathbf{x}_{1,1}, \dots, \mathbf{x}_{K,c_K}]$  to obtain the alignment-enhanced representations  $\mathbf{R} = [r_{1,1}, \dots, r_{K,c_K}]$  for the words in opinion clusters:

$$\mathbf{r}_{i,p} = \mathbf{W}_{rc} \text{aln}_{x_{i,p}} + \mathbf{x}_{i,p}, \quad (13)$$

where  $\mathbf{W}_{rc} \in \mathbb{R}^{d_e}$  is a model parameter. Figure 5 summarizes the construction.

At each decoding step  $q$  of the  $j$ -th opinion tag, the decoder reads the embeddings of the last prediction  $y_{j,q-1}$ .

$$y_{j,q} = \text{transformer\_decoder}(y_{j,q-1}), \quad (14)$$

where  $y_{j,q} \in \mathbb{R}^{d_e}$  is the target word representation. Next, we introduce our decoder in detail.

To capture semantic and alignment information from the opinion clusters, a multi-head cross-attention MHAtt [46] is applied to compute the attention score  $[\alpha_{1,1}, \dots, \alpha_{k,L_k}]$  between the last prediction  $y_{j,q-1}$  and  $[\mathbf{r}_{1,1}, \dots, \mathbf{r}_{k,c_k}]$ :

$$\mathbf{u}_{j,q-1}^z = \mathbf{W}_a^z y_{j,q-1} \quad (15)$$

$$\mathbf{k}_{i,p}^z = \mathbf{W}_b^z \mathbf{r}_{i,p} \quad (16)$$

$$\alpha_{i,p}^z = \frac{\exp(\mathbf{u}_{j,q-1}^z \cdot \mathbf{k}_{i,p}^z)}{\sum_{i'=1}^K \sum_{p'=1}^{c_{i'}} \exp(\mathbf{u}_{j,q-1}^z \cdot \mathbf{k}_{i',p'}^z)}, \quad (17)$$

where  $\mathbf{u}_{j,q-1}^z \in \mathbb{R}^{d_h}$ ,  $\mathbf{k}_{i,p}^z \in \mathbb{R}^{d_h}$  are query and key vectors that are linearly transformed from  $y_{j,q-1}$  and  $\mathbf{r}_{i,p}$  as in [46];  $z \in \{1, \dots, n_h\}$  indicates the  $z$ -th head among  $n_h$  heads;  $d_h = d_e/n_h$  is the dimension of each head.

The attention scores  $[\alpha_{i,j}^1, \dots, \alpha_{i,j}^{n_h}]$  are then used to compute an aggregated vector  $\mathbf{c}_{j,q}$  for target word  $y_{j,q}$ :

$$\mathbf{c}_{j,q} = \left[ \sum_{i=1}^K \sum_{j=1}^{c_i} \alpha_{i,j}^1 \mathbf{r}_{i,j}; \dots; \sum_{i=1}^K \sum_{j=1}^{c_i} \alpha_{i,j}^{n_h} \mathbf{r}_{i,j} \right]. \quad (18)$$

Then we feed the last word representation  $y_{j,q-1}$  and vector  $\mathbf{c}_{j,q}$  to a two-layer feed-forward network with a ReLU as the activation function and a highway layer normalization on top:

$$\mathbf{s}_{j,q-1} = \text{LayerNorm}(y_{j,q-1}^{n-1} + \mathbf{c}_{j,q}) \quad (19)$$

$$y_{j,q-1}^n = \text{LayerNorm}(\mathbf{s}_{j,q-1} + \text{FFN}(\mathbf{s}_{j,q-1})), \quad (20)$$

where  $y_{j,q} = y_{j,q-1}^n$  is the target word representation. After that, we use  $y_{j,q}$  to compute a probability distribution over the words in a predefined vocabulary  $\mathcal{V}$ , as shown in the following equation:

$$P_{\mathcal{V}}(y_{j,q} | [y_{1,0}, \dots, y_{j,q-1}], \mathcal{X}) = \text{softmax}(\mathbf{W}_v y_{j,q} + \mathbf{b}_v), \quad (21)$$

where  $\mathbf{W}_v \in \mathbb{R}^{|\mathcal{V}| \times d_e}$ ,  $\mathbf{b}_v \in \mathbb{R}^{|\mathcal{V}|}$  are trainable parameters. To enable our model to generate out-of-vocabulary (OOV) words, we adopt the copy mechanism [38] to predict OOV words by directly copying words from the opinion clusters. We first compute a soft gate  $p_{gen} \in [0, 1]$  between generating a word from the predefined vocabulary  $\mathcal{V}$  and copying a word from the input reviews  $\mathcal{X}$ :

$$p_{gen} = \sigma(\mathbf{W}_g y_{p,q} + b_g), \quad (22)$$

where  $\mathbf{W}_g \in \mathbb{R}^{d_e}$  and  $b_g \in \mathbb{R}$  are trainable parameters. Finally, we can derive the final next-word probability distribution  $P(y_{j,q})$ :

$$\alpha_{i,p} = \frac{\sum_z^{n_h} \alpha_{i,p}^z}{n_h}, \quad (23)$$

$$P(y_{j,q}) = p_{gen} P_{\mathcal{V}}(y_{j,q}) + (1 - p_{gen}) \sum_{i,p:x_{i,p}=y_{j,q}} \alpha_{i,p}, \quad (24)$$

where we use  $P(y_{j,q})$  to denote  $P(y_{j,q} \mid [y_{1,0}, \dots, y_{j,q-1}], \mathcal{X})$  for brevity. We use the negative log-likelihood of the ground-truth words  $y_{j,q}^*$  as the generation loss function:

$$\mathcal{L}_{gen} = - \sum_{j=1}^N \sum_{q=1}^{L_{y_j}} \log P(y_{j,q}^* \mid [y_{1,0}^*, \dots, y_{j,q-1}^*], \mathcal{X}). \quad (25)$$

*Alignment Loss.* During the generation of the  $j$ -th opinion tag, we enforce an alignment loss to help locate FOCs accurately. The model is explicitly taught to focus on FOCs and ignore OOCs in the attention  $\alpha_{i,p}$  via the following loss:

$$\mathcal{L}_{aln} = - \log \left( \frac{\sum_{i,p:x_{i,p} \in \text{FOCs}} \alpha_{i,p}}{\sum_{i,p} \alpha_{i,p}} \right) + \log \left( \frac{\sum_{i,p:x_{i,p} \in \text{OOCs}} \alpha_{i,p}}{\sum_{i,p} \alpha_{i,p}} \right). \quad (26)$$

## 4.5 Multi-task Training Objective

We adopt a multi-task learning framework to jointly minimize the salience classification loss, alignment loss, and generation loss. The objective function is:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{cla} + \lambda_2 \mathcal{L}_{aln} + \lambda_3 \mathcal{L}_{gen}, \quad (27)$$

where  $\lambda_1, \lambda_2, \lambda_3$  are hyper-parameters that control the weights of these three losses. We set  $\lambda_1 = \lambda_2 = \lambda_3 = 1$ . Thus, each component of our joint model can be trained end-to-end.

## 5 EXPERIMENTAL SETUP

We set up experiments to compare AOT-Net against a number of relevant baselines. We are interested in the overall performance of AOT-Net and in understanding the effectiveness of the salience estimation and ranking alignment.

### 5.1 Experiments

We report on five experiments. First, we compare AOT-Net against a number of baselines to assess its overall performance. Then we conduct ablation studies to analyze the influence of different components in AOT-Net as follows: (i) **w/o SSE** is AOT-Net without the sentence-level salience estimating component (SSE). (ii) **w/o RCR** is AOT-Net without the review clustering and ranking component (RCR). (iii) **w/o AF** is AOT-Net without alignment feature (AF). (iv) **w/o AL** is AOT-Net without alignment loss (AL). To further explore the effectiveness of the sentence-level self-attention mechanism in SSE, we consider **AOT-RNN**, the method only considers BiGRU in salience score prediction; whereas we write **AOT-Embed** for the method that employs MLP to replace BiGRU in SSE. Fourth, we analyze the performance of AOT-Net for different sizes of FOCs. Lastly, we provide a case study about abstractive opinion tagging.

### 5.2 Baselines

We compare AOT-Net with the following methods: (i) **TF-IDF** is an extractive approach that selects the important words as summary based on term frequency and inverse document frequency; (ii) **TextRank** [35] is an unsupervised algorithm based on weighted-graphs; (iii) **RNN** is a sequence to sequence model with attention implemented by bi-directional GRU layer [12]; (iv) **PG-Net** [38] is a classical opinion summarization model based on the encoder-decoder framework with attention and copy mechanisms; and (v) the **Transformer** [46] is a Transformer-based encoder-decoder model with a copy mechanism, which is a strong baseline widely-adopt in opinion summarization.

### 5.3 The eComTag Dataset

Since there is no available opinion tagging dataset, we build a new one, named *eComTag* from several Chinese e-commerce websites. We collect nearly 112k items, sampling from different domains, including Cosmetic (37.43%), Electronics (29.51%), Books (10.57%), Entertainment (8.23%), Food (7.62%), Sports (3.96%), Clothes (3.16%), Medical (1.62%), and Furniture (0.33%). For each domain, there are a set of reviews and a list of opinion tags. Since reviewers may comment on multiple aspects, e.g., “The dim sum tasted extremely fresh, and the price was quite reasonable!”, we split each review into sentences by punctuation. We use a sentence to denote a review in our paper. Then, we remove samples where the number of opinion tags is smaller than 4 or the number of reviews is fewer than 50. Finally, we construct *eComTag* with 50,068 item samples. Users may write reviews arbitrarily, which results in many meaningless expressions, such as “Love, love, LOVE this space!”, and “come with my boyfriend.” To teach our model to distinguish these noisy sentences, we annotate each review with a binary salience label via human judgment. If a review is item-related, we label the review as 1, otherwise 0. The salience labels are the supervision signals for the sentence-level salience estimating component.

Finally, each item sample consists of a set of reviews, a set of corresponding salience labels, and a sequence of opinion tags. For text preprocessing, we tokenize texts using the Jieba toolkit<sup>6</sup> and maintain a 50k vocabulary. In *eComTag*, about 30% of samples have more than 1024 words in reviews. We randomly split the dataset into training/validation/test sets with 8:1:1 ratio. The statistics are shown in Table 1. Specially, we define the *present tag* (Pr) as the exact tag that appears in reviews and *absent tag* (Ab) as the tag unseen in reviews. The proportion of absent tags is close to 75%, which further proves the necessity to apply abstractive methods on opinion tagging.

### 5.4 Evaluation Metrics

We employ two information retrieval metrics to evaluate the opinion tag generation: the macro F@ $k$  score and the normalized discounted cumulative gain (NDCG@ $k$ ) score. Both are widely used to measure word overlap [9, 42]. To measure diversity of the generated opinion tags, we adopt the Distinct-2 score [22] and a macro Unique- $N$  score. We compute Unique- $N$  as follows, Unique- $N = \sum_{i=1}^T N_i / T$ , where  $N_i$  is the number of distinct opinion tags in the  $i$ -th sample.

<sup>6</sup><https://github.com/fxsjy/jieba>

**Table 1: Statistics of the eComTag dataset. “Pr” and “Ab” denote the proportion of present tags and absent tags respectively. “MTN” and “MTL” indicate max tag number and max tag length for a sample respectively.**

Data	Sample	Pr	Ab	MTN	MTL
Training	40,162	24.1%	75.9%	19	40
Validation	4,953	24.3%	75.7%	20	39
Test	4,953	24.1%	75.9%	14	32

Moreover, we design two metrics, Exact Rank Match (ERM) and Fuzzy Rank Match (FRM), to evaluate the rank accuracy of opinion tags. ERM is defined as the one-to-one exact match proportion between true tags and predicted tags. Inspired by the Embedding Score [27], FRM first maps predicted tags and the corresponding true tags into the same vector space, and then computes the average cosine similarity between their vector representations.

## 5.5 Implementation Details

We adopt the Adam [20] optimizer with settings  $\{\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, lr = 10^{-4}\}$  and we vary the learning rate following Vaswani et al. [46]. We add dropout [40] with keeping rate 0.8 and label smoothing [43] with smoothing factor 0.1. We use the Tencent AI Lab Chinese Embeddings<sup>7</sup> for initialization of the word embedding layers. The rest of the parameters are randomly initialized. The dimensions of the alignment feature, word embedding layers and positional embedding layers are set to 200. We set the batch size to 16 and use the validation loss for early stopping. When inference, we set the maximum decoding step as 50. We use a bidirectional GRU [12] with 2 layers to implement the sentence-level salience estimating component. All RNN-based models have 256 hidden units. All transformer-based models have 300 hidden units; the feed-forward hidden size is set to 50 for all layers. We set the  $F$  in Section 4.4 to 3 (3 is the number of *Focused Opinion Clusters* at each decoding step) to ensure the target tag token have enough relevant reviews to reference and avoid introducing too much interference information simultaneously. AOT-Net was trained on a single Tesla V100 GPU and is implemented using PyTorch. All hyperparameters and models are selected on the validation set and the results are reported on the test set.

## 6 RESULTS AND ANALYSIS

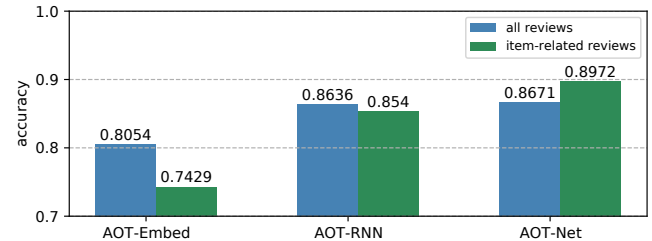
Overall evaluation results on generating opinion tags are listed in Table 2. We find that all the abstractive models significantly outperform all the traditional extractive baselines. Thus we conclude that informal and colloquial nature of user-generated reviews make item-related features indiscernible using the unsupervised extraction methods. As expected, we also find that PG-Net significantly outperforms RNN, which implies that the copying mechanism is useful for opinion summarization. We can see AOT-Net significantly outperforms baseline Transformer in terms of all metrics and achieves the best performance for most metrics. The results of our ablation studies are shown in the lower part of the Table 2. We observe that after removing the SSE component, performances of AOT-Net in terms of most metrics drops obviously. If we do

<sup>7</sup><https://ai.tencent.com/ailab/nlp/embedding.html>

not rank and group reviews into opinion clusters before decoding (i.e., w/o RCR), although the diversity metric has a slight increase, the rank accuracy of AOT-Net decreases as we anticipated. We also find that after removing alignment feature or alignment loss mechanisms in the decoder, the performance of both retrieval and rank accuracy metrics (i.e.,  $F_1$  and ERM) degrades. We will conduct a detailed analysis of the individual components in the following sections.

## 6.1 Salience Estimation Analysis

As shown in Table 2, AOT-Net achieves a 17.1% and 2.72% increase over “w/o SSE” in terms of Micro-Distinct-2 and Macro-Distinct-2, respectively. Similar improvements can be observed for other metrics. This demonstrates that the predicted salience scores help AOT-Net to focus on more valuable reviews. To verify the effectiveness of SSE with more details, in Figure 6 we list the accuracy scores of AOT-Embed, AOT-RNN, and AOT-Net for sentence-level salience estimation. We find that AOT-Net outperforms both AOT-Embed and AOT-RNN, which verifies the effectiveness of BiGRU and self-attention mechanisms. AOT-RNN achieves 7.2% increase over AOT-Embed in terms of accuracy for all reviews, which verifies the advantage of BiGRU in representing review. In terms of accuracy, we find that AOT-Net gives a 0.4% and 5.1% increase over AOT-RNN for all reviews and item-related reviews, respectively. This indicates that AOT-Net benefits from self-attention mechanisms, which capture the shared information to distinguish item-related reviews.



**Figure 6: Accuracy values for sentence-level salience estimation.**

## 6.2 Number of FOCs

To evaluate the effect of the number of FOCs on the performance of rank-aware opinion tagging, we examine the performance of AOT-Net with different values of  $F$  (see Section 4.4) in terms of ERM, FRM, and Distinct-2, respectively. As shown in Table 3,  $F = 1$  significantly decreases the model rank and diversity performance. This suggests that only focusing on a single opinion cluster ignores many related reviews. We also find that when  $F = 5$ , the performance of AOT-Net slightly decreases; AOT-Net achieves the best performance in terms of all metrics when  $F = 3$ . Hence, we infer that  $F$  is a trade-off between focusing on relevant reviews and removing irrelevant noise.

## 6.3 Case Study

Figure 7 shows an example illustrating the 5 highest attention weights  $\alpha_{i,p}$  during the generation of opinion tags. We see that the

**Table 2: Evaluation results on the *eComTag* dataset. Results in bold are leading results in terms of the corresponding metric.**

Models	F <sub>1</sub> @5	F <sub>1</sub> @10	NDCG@5	NDCG@10	ERM	FRM	Distinct-2		
							Micro	Macro	Unique-N
TF-IDF	0.0039	0.0038	0.0168	0.0169	0.16	0.19	–	–	–
TextRank	0.0019	0.0018	0.0091	0.0097	0.06	0.21	–	–	–
RNN	0.2895	0.2753	0.7383	0.7701	0.17	0.44	0.60	62.19	7.447
PG-Net	0.3138	0.2896	<b>0.7600</b>	0.8009	0.19	0.44	1.33	65.78	6.798
Transformer	0.2833	0.2756	0.6916	0.7483	0.25	0.59	1.57	89.23	8.851
<b>AOT-Net</b>	<b>0.3529</b>	<b>0.3492</b>	0.7473	<b>0.8045</b>	<b>0.31</b>	<b>0.64</b>	1.30	<b>94.35</b>	8.953
w/o SSE	0.2930	0.2822	0.7022	0.7563	0.25	0.61	1.11	91.85	8.957
w/o RCR	0.3434	0.3370	0.7353	0.7913	0.31	0.63	<b>1.77</b>	92.94	8.935
w/o AF	0.3141	0.3056	0.7194	0.7768	0.28	0.62	1.21	93.23	<b>8.997</b>
w/o AL	0.3406	0.3336	0.7322	0.7857	0.30	0.64	1.30	93.11	8.926

**Table 3: Performance on different numbers (*F*) of FOCs.**

Models	ERM	FRM	Macro	Distinct-2
<i>F</i> =1	0.30	0.62		93.24
<i>F</i> =3	0.31	0.64		94.35
<i>F</i> =5	0.30	0.63		93.32

Good	service	attitude	.	The	service	of	hairdresser
was	pretty	good	!	Very	satisfactory	service	.
In	a	word	:	the	service	is	very
good	!	The	service	is	exceptional	and	the
hairstyle	is	good	.	Good	service	from	boss
.	Many	tools	but	they	are	very	skilled
.	Hair	coloring	works	well	.	The	manager
is	very	skilled	.	The	result	is	satisfactory
.	The	hair	looks	good	.	The	result
is	what	I	expected	.	I'll	cut	my
hair	here	.	The	price	is	better	online
.	The	price	is	acceptable	.	The	staff
is	enthusiastic	,	I'm	satisfied	.	Nice	haircut
experience	!	We	had	a	great	experience	!

Prediction: 1<sup>st</sup> service enthusiasm; 2<sup>nd</sup> professional staffs; 3<sup>rd</sup> good effect  
4<sup>th</sup> very affordable; 5<sup>th</sup> great experience  
Reference: service enthusiasm, professional staffs, good effect, great experience, very affordable.

**Figure 7: The transition of reviews attention distribution between tags computed by AOT-Net. Different colors correspond to different tags. We only depict *top@5* attention probabilities for brevity. (Best viewed in color).**

model transits its focus smoothly from the first review sentence to later review sentences. Sometimes, the model may focus on the same review for two opinion tags, such as “The *service* (1st tag) is good and *hairstyle* (2nd tag) is good.” The first three predicted tags were completely accurate. However, the transformer only predicted the first tag accurately. To validate the effectiveness of the alignment loss, we calculate  $\sum_{i,p:x_{i,p} \in \text{FOCs}}$  and  $\sum_{i,p:x_{i,p} \in \text{OOCs}}$  for all items in the test set. Results show that  $\sum_{i,p:x_{i,p} \in \text{FOCs}}$  and  $\sum_{i,p:x_{i,p} \in \text{OOCs}}$  in AOT-Net are 0.7304 and 0.2696 on average. However, for AOT-Net w/o AL,  $\sum_{i,p:x_{i,p} \in \text{FOCs}}$  and  $\sum_{i,p:x_{i,p} \in \text{OOCs}}$  are 0.6509 and 0.3491, respectively. This phenomenon indicates that the alignment feature by itself is not enough to force AOT-Net to focus on FOCs. In particular, we compare the ranks of opinion tags generated by the transformer and AOT-Net respectively. For the transformer, items where the first 3 opinion tags have accurate ranks only account for near 6.11%, while for AOT-Net, the proportion can reach 9.68%.

Thus, we conclude that the alignment feature and alignment loss in AOT-Net are helpful to capture the ranks of the opinion tags.

## 7 CONCLUSION AND FUTURE WORK

In this paper, we have proposed the *abstractive opinion tagging* task, which aims to automatically generate a ranked list of opinion tags from a large number of reviews. We have proposed a rank-aware abstractive opinion tagging framework (AOT-Net) that includes a sentence-level salience estimating component, a review clustering and ranking component, and a rank-aware opinion tagging component. To validate the effectiveness of AOT-Net, we conduct extensive experiments on a newly collected real-world dataset, *eComTag*. Experiments show that AOT-Net achieves state-of-the-art performance on the abstractive opinion tagging task. AOT-Net has two main advantages over previous work. On the one hand, it generates more concise opinion tags; on the other hand, the ranked lists of generated opinion tags help users distinguish products with very similar aspects. Our work provides a plausible solution to greatly reduce human annotation costs for online e-commerce opinion tagging. Although we focused mostly on e-commerce portals, our methods are also broadly applicable to other settings with opinionated content, such as microblogs.

Limitations of our work include its low efficiency and coarse-grained salience estimation. As to our future work, we will adopt the deep clustering network instead of *k*-means algorithm. Also, pre-trained language models could provide more power to enhance our sentence salience estimation. Few-shot learning for handling unbalanced review distributions among different domains could be another direction. It will be also interesting to explore user interactions with AOT-Net to generate personalized opinion tags in the future.

## CODE AND DATA

The source code and dataset used in this paper are available at <https://github.com/qtli/AOT>.

## ACKNOWLEDGMENTS

We thank our reviewers for valuable feedback. This work was supported by the National Key R&D Program of China with grant No. 2020YFB1406704, the Natural Science Foundation of China



(61972234, 61902219, 61672324, 61672322, 62072279), the Key Scientific and Technological Innovation Program of Shandong Province (2019JZZY010129), the Tencent AI Lab Rhino-Bird Focused Research Program (JR201932), the Fundamental Research Funds of Shandong University, the Foundation of State Key Laboratory of Cognitive Intelligence, iFLYTEK, P.R. China (COGOSC-20190003). All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

## REFERENCES

- [1] Reinald Kim Amplayo and Mirella Lapata. 2019. Informative and Controllable Opinion Summarization. *CoRR abs/1909.02322* (2019).
- [2] Reinald Kim Amplayo and Mirella Lapata. 2020. Unsupervised Opinion Summarization with Noising and Denoising. In *ACL*. 1934–1945.
- [3] Stefanos Angelidis and Mirella Lapata. 2018. Summarizing Opinions: Aspect Extraction Meets Sentiment Prediction and They Are Both Weakly Supervised. In *EMNLP*. 3675–3686.
- [4] Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer Normalization. *CoRR abs/1607.06450* (2016).
- [5] Arthur Brazinskis, Mirella Lapata, and Ivan Titov. 2020. Unsupervised Opinion Summarization as Copycat-Review Generation. In *ACL*. 5151–5169.
- [6] Giuseppe Carenini, Jackie Chi Kit Cheung, and Adam Pauls. 2013. Multi-Document Summarization of Evaluative Text. *Comput. Intell.* 29, 4 (2013), 545–576.
- [7] Giuseppe Carenini, Raymond T. Ng, and Adam Pauls. 2006. Multi-Document Summarization of Evaluative Text. In *EACL*. The Association for Computer Linguistics.
- [8] Nofar Carmeli, Xiaolan Wang, Yoshihiko Suhara, Stefanos Angelidis, Yuliang Li, Jinfeng Li, and Wang-Chiew Tan. 2020. ExplainIt: Explainable Review Summarization with Opinion Causality Graphs. *CoRR abs/2006.00119* (2020).
- [9] Hou Pong Chan, Wang Chen, and Irwin King. 2020. A Unified Dual-view Model for Review Summarization and Sentiment Classification with Inconsistency Loss. In *SIGIR*. 1191–1200.
- [10] Hou Pong Chan, Wang Chen, Lu Wang, and Irwin King. 2019. Neural Keyphrase Generation via Reinforcement Learning with Adaptive Rewards. In *ACL*. 2163–2174.
- [11] Wang Chen, Yifan Gao, Jiani Zhang, Irwin King, and Michael R. Lyu. 2019. Title-Guided Encoding for Keyphrase Generation. In *AAAI*, Vol. 33. 6268–6275.
- [12] Kyunghyun Cho, Bart van Merriënboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In *EMNLP*. 1724–1734.
- [13] Guy Elad, Ido Guy, Slava Novgorodov, Benny Kimelfeld, and Kira Radinsky. 2019. Learning to Generate Personalized Product Descriptions. In *CIKM*. 389–398.
- [14] Giuseppe Di Fabbri, Amanda Stent, and Robert J. Gaizauskas. 2014. A Hybrid Approach to Multi-document Summarization of Opinions in Reviews. In *INLG*. 54–63.
- [15] Kavita Ganesan, ChengXiang Zhai, and Jiawei Han. 2010. Opinions: A Graph Based Approach to Abstractive Summarization of Highly Redundant Opinions. In *COLING*. 340–348.
- [16] Shen Gao, Zhaochun Ren, Yihong Eric Zhao, Dongyan Zhao, Dawei Yin, and Rui Yan. 2019. Product-Aware Answer Generation in E-Commerce Question-Answering. In *WSDM*. 429–437.
- [17] Shima Gerani, Yashar Mehdad, Giuseppe Carenini, Raymond T. Ng, and Bitia Nejat. 2014. Abstractive Summarization of Product Reviews Using Discourse Structure. In *EMNLP*. 1602–1613.
- [18] Sujatha Das Gollapalli, Xiao-Li Li, and Peng Yang. 2017. Incorporating Expert Knowledge into Keyphrase Extraction. In *AAAI*.
- [19] Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *SIGKDD*. ACM, 168–177.
- [20] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*.
- [21] Tho Thi Ngoc Le, Minh Le Nguyen, and Akira Shimazu. 2016. Unsupervised Keyphrase Extraction: Introducing New Kinds of Words to Keyphrases. In *AI 2016 (LNCS)*, Vol. 9992. 665–671.
- [22] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A Diversity-Promoting Objective Function for Neural Conversation Models. In *NAACL*. 110–119.
- [23] Junjie Li, Xuepeng Wang, Dawei Yin, and Chengqing Zong. 2019. Attribute-aware Sequence Network for Review Summarization. In *EMNLP-IJCNLP*. 2998–3008.
- [24] Pengyuan Li, Lei Huang, and Guang-jie Ren. 2020. Topic Detection and Summarization of User Reviews. *CoRR abs/2006.00148* (2020).
- [25] Piji Li, Zihao Wang, Lidong Bing, and Wai Lam. 2019. Persona-Aware Tips Generation. In *The World Wide Web Conference*. 1006–1016.
- [26] Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam. 2017. Neural Rating Regression with Abstractive Tips Generation for Recommendation. In *SIGIR*. 345–354.
- [27] Chia-Wei Liu, Ryan Lowe, Iulian Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation. In *EMNLP*. 2122–2132.
- [28] Dayiheng Liu, Yeyun Gong, Jie Fu, Wei Liu, Yu Yan, Bo Shao, Daxin Jiang, Jiancheng Lv, and Nan Duan. 2020. Diverse, Controllable, and Keyphrase-Aware: A Corpus and Method for News Multi-Headline Generation. *arXiv preprint arXiv:2004.03875* (2020).
- [29] Rui Liu, Zheng Lin, and Weiping Wang. 2020. Keyphrase Prediction With Pre-trained Language Model. *arXiv preprint arXiv:2004.10462* (2020).
- [30] Yue Lu, ChengXiang Zhai, and Neel Sundaresan. 2009. Rated Aspect Summarization of Short Comments. In *WWW*. 131–140.
- [31] Yi Luan, Mari Ostendorf, and Hannaneh Hajishirzi. 2017. Scientific Information Extraction with Semi-supervised Neural Tagging. In *EMNLP*.
- [32] James MacQueen. 1967. Some Methods for Classification and Analysis of Multivariate Observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, Vol. 1. Oakland, CA, USA, 281–297.
- [33] Olena Medelyan, Eibe Frank, and Ian H Witten. 2009. Human-competitive Tagging using Automatic Keyphrase Extraction. In *EMNLP*. 1318–1327.
- [34] Rui Meng, Sanqiang Zhao, Shuguang Han, Daqing He, Peter Brusilovsky, and Yu Chi. 2017. Deep Keyphrase Generation. *CoRR abs/1704.06879* (2017).
- [35] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing order into text. In *EMNLP*. 404–411.
- [36] Slava Novgorodov, Ido Guy, Guy Elad, and Kira Radinsky. 2019. Generating Product Descriptions from User Reviews. In *WWW*. 1354–1364.
- [37] Jishnu Ray Chowdhury, Cornelia Caragea, and Doina Caragea. 2019. Keyphrase extraction from disaster-related tweets. In *WWW*. 1555–1566.
- [38] Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get To The Point: Summarization with Pointer-Generator Networks. In *ACL*. 1073–1083.
- [39] Dinghan Shen, Guoyin Wang, Wenlin Wang, Martin Renqiang Min, Qinliang Su, Yizhe Zhang, Chunyuan Li, Ricardo Henao, and Lawrence Carin. 2018. Baseline Needs More Love: On Simple Word-Embedding-Based Models and Associated Pooling Mechanisms. In *ACL*. 440–450.
- [40] Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *J. Mach. Learn. Res.* 15, 1 (2014), 1929–1958.
- [41] Yoshihiko Suhara, Xiaolan Wang, Stefanos Angelidis, and Wang-Chiew Tan. 2020. OpinionDigest: A Simple Framework for Opinion Summarization. In *ACL*. 5789–5798.
- [42] Zhiqing Sun, Jian Tang, Pan Du, Zhi-Hong Deng, and Jian-Yun Nie. 2019. DivGraphPointer: A Graph Pointer Network for Extracting Diverse Keyphrases. In *SIGIR*. 755–764.
- [43] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *CVPR*. 2818–2826.
- [44] Duyu Tang, Bing Qin, Ting Liu, and Yuekui Yang. 2015. User Modeling with Neural Network for Review Rating Prediction. In *IJCAI*.
- [45] Wenyi Tay. 2019. Not All Reviews Are Equal: Towards Addressing Reviewer Biases for Opinion Summarization. In *ACL*. 34–42.
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*. 6000–6010.
- [47] Lu Wang and Wang Ling. 2016. Neural Network-Based Abstract Generation for Opinions and Arguments. In *NAACL*. 47–57.
- [48] Minmei Wang, Bo Zhao, and Yihua Huang. 2016. PTR: Phrase-based Topical Ranking for Automatic Keyphrase Extraction in Scientific Publications. In *ICONIP*. 120–128.
- [49] Yue Wang, Jing Li, Hou Pong Chan, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. Topic-Aware Neural Keyphrase Generation for Social Media Language. In *ACL*. 2516–2526.
- [50] Yue Wang, Jing Li, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. Microblog Hashtag Generation via Encoding Conversation Contexts. In *NAACL-HLT*. 1624–1633.
- [51] Wenting Xiong and Diane J. Litman. 2014. Empirical Analysis of Exploiting Review Helpfulness for Extractive Summarization of Online Reviews. In *COLING*. 1985–1995.
- [52] Qi Zhang, Yang Wang, Yeyun Gong, and Xuanjing Huang. 2016. Keyphrase Extraction Using Deep Recurrent Neural Networks on Twitter. In *EMNLP*. 836–845.
- [53] Yingyi Zhang, Jing Li, Yan Song, and Chengzhi Zhang. 2018. Encoding Conversation Context for Neural Keyphrase Extraction from Microblog Posts. In *NAACL-HLT*. 1676–1686.