# Understand me, if you refer to Aspect Knowledge: Knowledge-aware Gated Recurrent Memory Network

Bowen Xing and Ivor W. Tsang, *Fellow, IEEE*

*Abstract*—Aspect-level sentiment classification (ASC) aims to predict the fine-grained sentiment polarity towards a given aspect mentioned in a review. Despite recent advances in ASC, enabling machines to preciously infer aspect sentiments is still challenging. This paper tackles two challenges in ASC: (1) due to lack of aspect knowledge, aspect representation derived in prior works is inadequate to represent aspect's exact meaning and property information; (2) prior works only capture either local syntactic information or global relational information, thus missing either one of them leads to insufficient syntactic information. To tackle these challenges, we propose a novel ASC model which not only end-to-end embeds and leverages aspect knowledge but also marries the two kinds of syntactic information and lets them compensate for each other. Our model includes four key components: (1) a knowledge-aware gated recurrent memory network recurrently integrates dynamically summarized aspect knowledge; (2) a dual syntax graph network combines both kinds of syntactic information to comprehensively capture sufficient syntactic information; (3) a knowledge integrating gate re-enhances the final representation with further needed aspect knowledge; (4) an aspect-to-context attention mechanism aggregates the aspect-related semantics from all hidden states into the final representation. Experimental results on several benchmark datasets demonstrate the effectiveness of our model, which overpass previous state-of-the-art models by large margins in terms of both Accuracy and Macro-F1. To facilitate further research in the community, we have released our source code at https://github.com/XingBowen714/KaGRMN-DSG_ABSA.

*Index Terms*—Sentiment Analysis, Entity Knowledge, Aspect Level, Memory Network.

## I. INTRODUCTION

Aspect-level sentiment classification (ASC) [1] is a fine-grained task of sentiment classification or emotion recognition [2]–[5]. ASC aims to infer the fine-grained sentiment of a given aspect mentioned in a review. Generally, an aspect is a noun phrase included in a review sentence. For example, in a review "It took so long to get the check, while the dinner is great.", there are two aspects (*check* and *dinner*) of opposite sentiments. ASC has received increasing attention and interest from both academia and the industry due to its

Bowen Xing is with Australian Artificial Intelligence Institute (AAII), University of Technology Sydney. Ivor W. Tsang is with the A*STAR Centre for Frontier AI Research (CFAR), and also with Australian Artificial Intelligence Institute (AAII), University of Technology Sydney. E-mail: Bowen.Xing@student.uts.edu.au, ivor_tsang@ihpc.a-star.edu.sg

wide applications in real-life scenarios such as dialog systems [6], online reviews [7] and social networks [8].

Prior works have noticed the importance of aspect-context interaction. Different kinds of attention mechanisms [9]–[11] are proposed to extract aspect-relevant semantics from the hidden states of context words. And more recently, syntactic information is widely leveraged to facilitate the interactions between the aspect and its related words that are distant in context sequence. Graph Convolutional Networks (GCN) [12], [13] and Graph Attention Networks (GAT) [14], [15] are adopted to encode the syntax graphs predicted by off-the-shelf dependency parsers. [12], [16] employed GCNs to capture the local syntactic information. [15] proposed relational multi-head attention (Relational MHA) to capture the global relational information between aspect and each context word.

However, little attention has been spent on aspect representation and its conveyed semantics. Aspect representation and its semantics not only guide the aspect-context interaction but also provide important clues for ASC. Despite its importance, in previous works [12], [15], [17], aspect representation is simply derived by pooling the hidden states of aspect words. In Sec. IV-H we empirically study the aspect representation generated by BERT [18], and two cases are shown in Table VIII. We can find that BERT cannot capture the exact meanings and property information of *Mountain Lion OS* and *iTune*, although it is one of the strongest language models. Merely relying on pre-trained large language models cannot obtain sufficiently effective and informative aspect representation, making it hard for machines to address ASC. In contrast, humans can easily handle ASC and we conjecture the key to master this task is to leverage the adequate aspect knowledge they often refer to as the clue. Thinking of and leveraging the aspect knowledge are instinctive reactions of humans when they read an aspect in a review. For example, there is a review "Just a not bad restaurant, because the cheese and chips are both very soft." With the knowledge of 'cheese' and 'chips', humans are aware that the former should be soft and the latter should be click (not soft). Hence it is easy for humans to infer the positive sentiment of 'cheese' and the negative sentiment of 'chips'. However, in contrast, in ASC models there is no such mechanism, and aspect knowledge has not been explored or leveraged. Inheriting this deficiency, the aspect representation and semantics derived by prior models may lose important aspect information, which hinders aspect sentiment reasoning and make ASC challenging for machines.

On the other hand, both GCN and Relational MHA are useful for modeling distinct syntax graphs, but they have respective

shortages: GCN is hard to capture the global relations between aspect and its non-adjacent context words on the syntax graph; Relational MHA fails to capture the local syntactic information among context words because they are isolated from each other on the star-shaped aspect-oriented syntax graph. However, prior works only consider one of them, resulting in insufficient syntactic information.

To tackle the aforementioned two challenges, we suggest that (1) aspect knowledge should be explicitly leveraged in ASC models; (2) both kinds of syntactic information should be combined to capture sufficient syntactic information. We observe that there is plenty of entity descriptions in popular and easily accessible knowledge bases, such as DBpedia[1] and Wikipedia[2]. From their statistics, there are about 6.6 Million and 50 Million entities in current DBpedia and Wikipedia datasets. These descriptions can sufficiently represent the entities' meanings and conveying a wealth of entities' knowledge. In ASC, aspects are always entities, making it more convenient to retrieve their descriptions.

In this work, we propose a **K**nowledge-**a**ware **G**ated **R**ecurrent **M**emory **N**etwork with **D**ual **S**yntax **G**raph Modeling (**KaGRMN-DSG**) model as our solution to the two challenges. Specifically, its novelty lies in three core modules. The **first** one is Knowledge-aware Gated Recurrent Memory Network (KaGR-MN) which recurrently integrates the aspect knowledge into aspect representation and then context memories. An aspect-to-description attention mechanism is devised to dynamically summarize the needed aspect knowledge from the aspect description regarding the current semantic state. An adaptive knowledge integrating gate is designed to adaptively integrate the summarized knowledge into aspect representation. Then a self multi-head attention is employed to contextualize the integrated knowledge and update the context memory bank. The **second** one is Dual Syntax Graph Network (DSG-Net), which marries the proposed Position-aware GCN and Relational MHA, then learns the dual syntactic interaction to comprehensively capture sufficient syntactic information. The **third** one is the knowledge integrating gate (KI Gate) which re-enhances the final representation with further needed knowledge.

We highlight our contributions as follows:
(1) Based on plenty of informative entity descriptions from easily accessible knowledge bases, we end-to-end embed and leverage the aspect knowledge to address ASC.
(2) We propose a novel KaGR-MN, which combines the advantages of LSTM, Transformer, and Memory Networks. It recurrently embeds and integrates beneficial aspect knowledge into aspect representation and all context memories.
(3) We propose a dual syntax graph network, in which the local syntactic information and global relational information are combined to comprehensively capture sufficient syntactic information.
(4) We conduct extensive experiments on three benchmark datasets. Results show that our model achieves new state-of-the-art performances, significantly outperforming previous best results. Ablation study and further analysis validate the effectiveness of our model.

## II. RELATED WORKS

In early studies [19], [20], sentiment classifiers were built by traditional machine learning algorithms which demanded labor-intensive feature engineering. Most recently proposed ASC models are based on neural networks which can automatically learn representations. Conventionally, neural ASC model contains an aspect encoder, a context encoder and an aspect-to-context attention mechanism [9]–[11], [17], [21].

Different kinds of networks are adopted as the encoder. As for LSTM, [22] employed two separated LSTMs to encode the aspect-left and -right word sequences and then combined the two last hidden sates for classification;

[23] proposed to leverage external document sentiment analysis corpus in a multi-task framework to enhance the context modeling of LSTM. Besides, convolutional neural networks (CNN) and Memory Networks (MNs) are also exploited as encoders. [24] introduced parameterized filters and parameterized gates into CNN to integrate aspect information for context encoding. [25] designed a gated CNN layer to extract the aspect-specific features from the context hidden states. Based on standard MNs, [26] proposed the target-sensitive memory networks to focus on the impact of aspect semantics on the classification.

The attention mechanism is utilized to extract aspect-related sentiment features via assigning a weight to each context word regarding its relevance to the aspect. [9] proposed an aspect-to-context attention and a context-to-aspect attention to study the interactions between the aspect and context. [11] proposed an algorithm to automatically mine useful supervised information for the attention mechanism through the training process.

However, attention mechanisms may hardly capture the important words which is far from the aspect in the input context. As the development of graph neural networks [27]–[30], recent works utilize graph convolution network (GCN) [12], [13], [16], [28] and graph attention network (GAT) [15], [27], [31] to model the syntax graph for shortening the distance between the aspect and its sentiment trigger words and leveraging the syntactical information. [12] employed LSTM as encoder and exploit GCN to capture local syntactic information via encoding the syntax graph produced by off-the-shelf dependency parsers. [15] proposed Relational MHA, which can capture the global dependency between the aspect and each context word via operating on the star-shaped aspect-oriented syntax graph.

To enhance the context modeling, [23] and [32] trained their models on both document-level sentiment classification and ASC tasks in the multi-task framework with a shared encoder. [33] proposed an aspect-aware LSTM which introduces aspect information into LSTM cells to generate better context hidden states in which more aspect-related information is retained and aspect-irrelevant information is discarded. As BERT has proven its power of language modeling on hetergenuous NLP tasks, more recently proposed work [13], [15] adopted BERT as the context encoder to obtain high-quality hidden states.

However, prior models neglect to leverage aspect knowledge, resulting in inadequate aspect semantics. And the syntactic

---

[1] https://wiki.dbpedia.org/
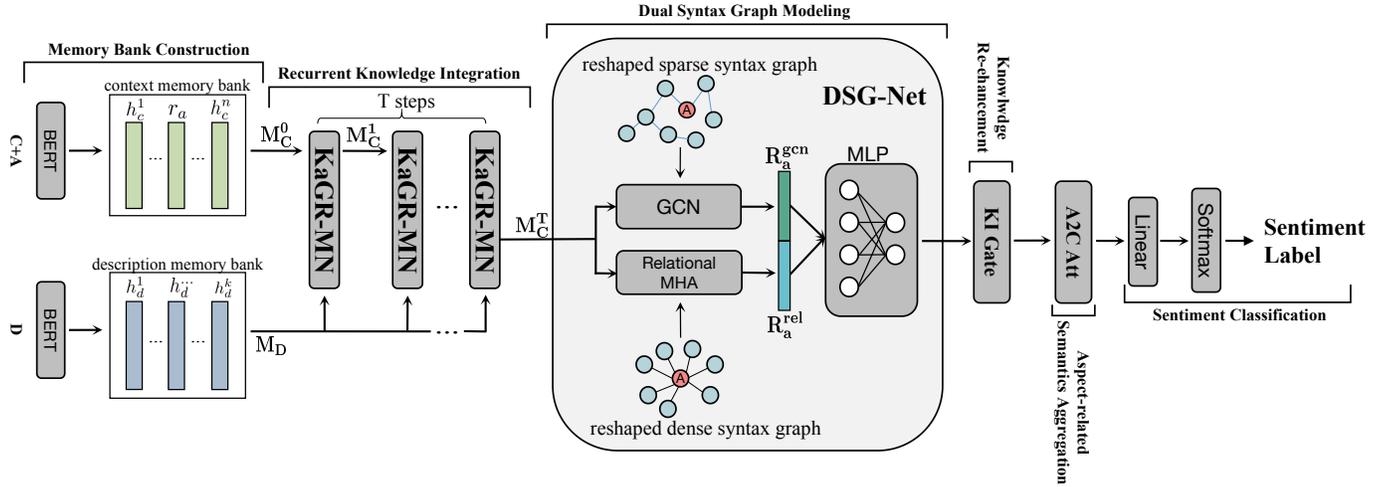
[2] https://www.wikipedia.org/

Fig. 1. The architecture of KaGRMN-DSG. The internal architecture of KaGR-MN cell is shown in Fig.2

information they captured is insufficient. In this paper, we propose KaGRMN-DSG to solve these two challenges. There are two main differences from our model and previous works. The first one is leveraging aspect knowledge, which is achieved by a novel KaGR-MN. The other one is combining both of local syntactic information and global relational information, which is achieved by a DSG-Net.

## III. KAGRMN-DSG

**Overview** The architecture of our KaGRMN-DSG model is illustrated in Fig. 1. To extract the beneficial clues for aspects, knowledge-aware gated recurrent memory network and knowledge integration gate incorporate summarized aspect knowledge to enrich aspect representation and all context memories. To capture sufficient syntactic information, dual syntax graph network combines local syntactic information and global relational information, then learns their mutual interaction. To comprehensively abstract high-level clues, aspect-to-context attention mechanism aggregates aspect-related semantics from all hidden states into the final representation. And we believe that these modules can effective cooperate to further improve aspect sentiment reasoning.

**Description Retrieval** We use aspect (**A**) to query DBpedia first and then Wikipedia to get its description (**D**). If multiple descriptions are returned (polysemy), the one with the highest semantic similarity to context (**C**) is selected as **D**. The semantic similarity of a description candidate and review context is calculated as:

$$avg(C) = \frac{1}{N_C} \sum_{i=1}^{N_C} e(c_i) \tag{1}$$

$$avg(D') = \frac{1}{N_D} \sum_{i=1}^{N_D} e(d_i) \tag{2}$$

$$sim(C, D') = cos\big(\alpha * avg(C) + (1 - \alpha) * e(dl), avg(D')\big) \tag{3}$$

where $N_C$ and $N_D'$ denotes the number of words in the context and description candidate respectively, $e(w)$ denotes

the word embedding[3] of word $w$, $dl$ denotes domain label (e.g. the $dl$ of Lap14 dataset is 'laptop'). Here we intuitively set $\alpha$ as 0.5 because both of the context semantics $avg(C)$ and domain information $e(dl)$ are important in selecting the correct description candidate. The reason why we use domain label here is that sometimes there may be not enough words conveying domain-specific semantics for distinguishing the needed description. Besides, the retrieval is enhanced with some rules, such as soft matching with lemmatization and stop word filtering. Finally, about 70% aspects in the datasets can be equipped with retrieved descriptions.

### A. Memory Bank Construction

In this work, we adopt BERT to encode the description and context to produce their hidden states. For description (**D**), the formal input is $\langle[\text{CLS}]; \mathbf{D}; [\text{SEP}]\rangle$, where $\langle;\rangle$ denotes concatenation operation. The description is encoded in the single-sentence manner then a series of its hidden states is generated: $\mathbf{H_D} = \{h_d^i \in \mathbb{R}^{d_e}\}_{i=1}^{N_D}$, which is taken as the description memory bank $\mathbf{M_D}$.

As for context (**C**), we model the context-aspect pair in the sentence-pair manner to generate aspect-aware hidden states [33]. The formal input is $\langle[\text{CLS}]; \mathbf{C}; [\text{SEP}]; \mathbf{A}; [\text{SEP}]\rangle$. In this way, we obtain the hidden state of $[\text{CLS}]$: $\mathbf{h_{cls}}$ and a series of aspect-aware context hidden states: $\mathbf{H_C} = \{h_c^i \in \mathbb{R}^{d_e}\}_{i=1}^{N_C}$. As BERT has a strong capability of sentence-pair modeling, $\mathbf{h_{cls}}$ contains not only the information from both of the aspect and the context but also their dependencies. Thus we take $\mathbf{h_{cls}}$ as the initial contextualized aspect representation $\mathbf{r_a^0}$. Then we use $\mathbf{r_a^0}$ to replace the hidden states of aspect words ($\mathbf{H_A} = \{h_a^i \in \mathbb{R}^{d_e}\}_{i=1}^{N_A}$) in $\mathbf{H_C}$, obtaining the initial context memory bank $\mathbf{M_C^0} = [h_c^1, h_c^2, ..., \mathbf{r_a^0}, ..., h_c^N]$, where $N = N_C - N_A + 1$.

$\mathbf{M_D}$ and $\mathbf{M_C^0}$ are two strands of input of KaGR-MN cell. Along time steps, $\mathbf{M_C}$ is recurrently updated while $\mathbf{M_D}$ remains identical.

---

[3]We use Glove word embedding [34].

## B. Knowledge-aware Gated Recurrent Memory Network

As the series of context hidden states and description hidden states have been obtained, now the challenge is *how to incorporate as much beneficial aspect knowledge as possible without losing the original semantics obtained from BERT?*.

The first thing is to conserve the original semantics in the context memories obtained from BERT. To this end, we employ Memory Networks (MNs) as the backbone to store context memories, because that MNs can accurately remember original facts [35]. Secondly, we are supposed to make the integrated knowledge beneficial. In other words, we should provide each sample the aspect knowledge it needs. Hence we propose an aspect-to-description attention (A2D Att) mechanism to summarize the needed aspect knowledge from the description memory bank. Thirdly, we should integrate the beneficial aspect knowledge into the aspect representation. Then we propose an adaptive knowledge integration gate, which borrows the idea of gating mechanisms in LSTM [36]. Gate mechanism has proven its strong ability of information integration in many tasks [33], [37]. However, only integrate knowledge into aspect representation is insufficient, not exploring the full value of aspect knowledge. It is intuitive that the aspect knowledge should be incorporated into all context memories. Besides, an appropriate mechanism should be devised to update the context memory bank. To achieve these two goals, inspired by Transformer [38], we utilize self multi-head attention to update the context memories, and in the meanwhile, the aspect knowledge in the aspect representation can be spread to all context memories. Finally, all the above mechanisms form the Knowledge-aware Gated Recurrent Memory Network (KaGR-MN), which combines the advantages of MNs, LSTM, and Transformer.

The architecture of KaGR-MN cell is illustrated in Fig. 2. In the following texts, we depict the details of KaGR-MN.

*1) Dynamic Knowledge Summarizing:* Intuitively, on the one hand, the context-aspect pair of each sample may demand individual aspect knowledge, even if they have the same aspect. On the other hand, at each time step, KaGR-MN should integrate specifically needed aspect knowledge according to the current cell state. Therefore, the aspect knowledge summarizing should be dynamic. To achieve this, we design an aspect-to-description attention (A2D Att) mechanism to dynamically summarize the specifically needed aspect knowledge from the description memory bank $\mathbf{M_D}$ at each time step. The architecture of A2D Att is shown in Fig. 2. At each time step ($t$), the aspect representation of previous time step $\mathbf{r_a^{t-1}}$ serves as the cell state and is used to query $\mathbf{M_D}$. Then an attention weight $\alpha$ is assigned to each $h_d$ regarding its importance to $\mathbf{r_a^{t-1}}$:

$$
\begin{aligned}
\alpha_i &= \texttt{SoftMax}(\mathcal{F}(h_d^i, r_a^{t-1})) \\
&= \frac{\exp(\mathcal{F}(h_d^i, r_a^{t-1}))}{\sum_{k=1}^{N_D} \exp(\mathcal{F}(h_d^k, r_a^{t-1}))}
\end{aligned}
\tag{4}
$$

where $\mathcal{F}(h_d^i, r_a^{t-1}))$ is a score function defined as:

$$
\mathcal{F}(h_d^i, r_a^{t-1})) = (\mathbf{W_d}\, h_d^i + \mathbf{b_d})\,(\mathbf{r_a^{t-1}})^{\mathsf{T}}
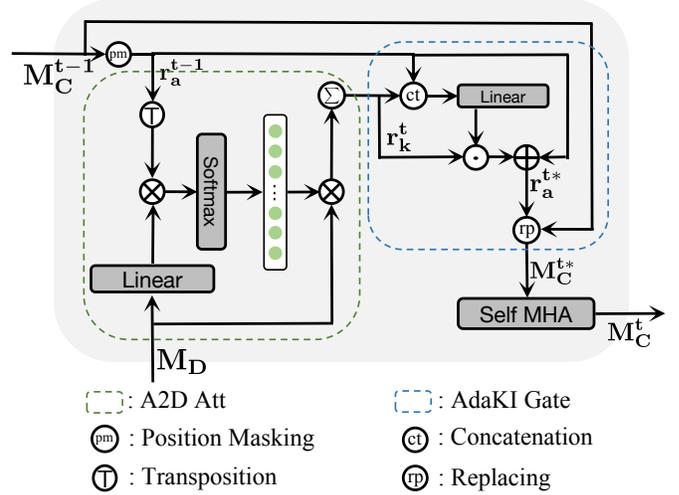\tag{5}
$$



Fig. 2. The architecture of KaGR-MN cell.

where $\mathbf{W_d}$ and $\mathbf{b_d}$ are weight matrix and bias respectively, and $\mathsf{T}$ denotes transposition. Then we can obtain the summarized knowledge representation as: $\mathbf{r_k^t} = \sum_{i=1}^{N_D} \alpha_i h_d^i$.

*2) Adaptive Knowledge Integration:* As the specifically needed knowledge has been summarized, it should be integrated into the aspect representation regarding the current cell state. As the gate mechanisms [36], [37], [39] have proven their ability of controlling information flow, here we design an **Ada**ptive **K**nowledge **I**ntegration (AdaKI) Gate to integrate $\mathbf{r_k^t}$ into $\mathbf{r_a^{t-1}}$. Its architecture is shown in Fig. 2. AdaKI Gate can be formulated as:

$$
\mathbf{r_a^{t*}} = \mathbf{r_a^{t-1}} + \mathbf{r_k^t} \odot (\mathbf{W_k}[\mathbf{r_a^{t-1}}, \mathbf{r_k^t}])
\tag{6}
$$

where $\odot$ denotes Hadamard product, $[,]$ denotes concatenation and $\mathbf{W_k}$ is weight matrix. The core of AdaKI Gate is to produce a gate vector using $\mathbf{r_k^t}$ and $\mathbf{r_a^{t-1}}$. This gate vector achieves the fine-grained control on each dimension of $\mathbf{r_k^t}$.

There are two merits of this fine-grained control. First, AdaKI Gate can determine what knowledge and how much knowledge from $\mathbf{r_k^t}$ should be integrated into $\mathbf{r_a^{t-1}}$. Second, it can map the integrated knowledge into the same semantic space of $\mathbf{r_a^{t-1}}$ and $\mathbf{M_C^{t-1}}$. This adaption helps maintain the semantics consistency of $\mathbf{r_a^{t*}}$ and $\mathbf{M_C^{t-1}}$, which is beneficial to later knowledge contextualizing. In Sec. IV-F, we investigate the effect of different knowledge gates used here. After $\mathbf{r_a^{t*}}$ is obtained, it replaces $\mathbf{r_a^{t-1}}$ in $\mathbf{M_C^{t-1}}$, forming $\mathbf{M_C^{t*}}$.

*3) Knowledge Contextualizing and Context Memory Bank Updating:* Although the needed beneficial knowledge has been integrated into $\mathbf{r_a^{t*}}$, the other context memories in $\mathbf{M_C^{t*}}$ remain the same as the ones in $\mathbf{M_C^{t-1}}$. Intuitively, all context memories should benefit from aspect knowledge to facilitate aspect-related information aggregation. To achieve this, we propose a knowledge contextualizing mechanism to broadcast the newly-integrated knowledge in $\mathbf{r_a^{t*}}$ to all context memories in $\mathbf{M_C^{t*}}$. Here we borrow the idea of self-attention [38], [40], which can effectively relate the different tokens in a sentence and capture the intra-sentence dependencies.

In this work, we adopt the self multi-head attention (Self MHA) formulation in [38]. We first map $\mathbf{M_C^{t*}}$ to queries

($\mathbf{Q}$), keys ($\mathbf{K}$) and values ($\mathbf{V}$) matrices by individual linear projections, where $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d_s}$. And this process repeats $H_n^s$ times, where $H_n^s$ is the number of heads and $H_n^s \times d_s = d_e$. The scaled dot-product attention is used to produce the output of each head, then all of the $H_n^s$ outputs are concatenated to form the updated context memory bank $\mathbf{M_C^t}$:

$$\mathbf{M_C^t} = \|_{h=1}^{H_n^s} \mathtt{SoftMax}\left(\frac{\mathbf{QK^T}}{\sqrt{d_s}}\right)\mathbf{V} \qquad (7)$$

This simple but effective knowledge contextualizing mechanism updates $\mathbf{M_C^{t*}}$ and $\mathbf{r_a^{t*}}$ by letting context memories (including $\mathbf{r_a^{t*}}$) exchange useful information with each other, which is beneficial to capture aspect-related information. Along time steps, $\mathbf{r_a}$ and $\mathbf{M_C}$ would contain more and more reasonable and beneficial semantics for ASC.

### C. Dual Syntax Graph Network

As proven in prior works, GCN can capture local syntactic information and Relational MHA can capture the global relation between the aspect and each context word via operating on the star-shaped aspect-oriented syntax graph (as shown in Fig. 1). However, as we have discussed in Sec. I, they have respective shortages. They can only capture one of the two kinds of syntactic information and lose the other. Previous models only employ one of them, leading to insufficient syntactical information.

To this end, we propose DSG-Net (as shown in Fig. 1) which marries the proposed Position-aware GCN and Relational MHA and learns their interaction, capturing sufficient syntactic information.

*1) Local Syntactic Information Modeling:* **Graph Construction** Based on the original syntax graph $G^4$, we first add a new aspect node $A$ and merge all edges between nodes of aspect words and non-aspect context words to $A$. Then we delete all of the original nodes of aspect words and their edges. The obtained graph is similar to $G$, and only several context word nodes are connected to $A$. Thus we term it sparse graph $G_s$ (shown in Fig.1).
**Position-aware GCN** In this work, we augment the standard GCN with a position weight $w_p^i = 1 - \frac{|i - \tau|}{N+1}$, in which $\tau$ denotes the position of aspect, $i$ denotes the $i^{th}$ context word. As the Self MHA in KaGR-MN does not consider the order of context memories, some positional and ordering information may be lost. $w_p^i$ can supplement this information, which helps capture local syntactic information. Besides, it indicates the position of $A$ and highlights the potential aspect-related words which are generally closer to $A$. In $l^{th}$-layer, the local neighborhood information is aggregated as:

$$h_i^l = \sum_{j \in \mathcal{N}_i^s} \mathbf{W_g^s}(w_p^j \, h_j^{l-1})/(d_i + 1) + \mathbf{b_g^s} \qquad (8)$$

in which $\mathcal{N}_i^s$ is the first-order neighbors of node $i$ (including $i$) in $G_s$, $d_i$ is the degree of node $i$, $\mathbf{W_g^s}$ and $\mathbf{b_g^s}$ are weight matrix and bias.

---

[4]obtained by spaCy toolkit: https://spacy.io/

*2) Global Relational Information Modeling:* We obtain the star-shaped aspect-oriented syntax graph following [15]. In this syntax graph, every context word directly connects to the aspect node $A$, so we term it dense graph $G_d$. Then we employ the Relational MHA to model the global relational dependency between aspect and each context word. The node representation is:

$$h_i = \sum_{m=1}^{H_n^d}\left(\sum_{j \in \mathcal{N}_i^d} \beta_{ij}^m \mathbf{W_m^1} h_j\right)/H_n^d$$
$$\beta_{ij}^m = \mathtt{SoftMax}(g_{ij}^m) \qquad (9)$$
$$g_{ij}^m = \mathtt{ReLU}\left(r_{ij}\mathbf{W_m^2} + \mathbf{b_m^1}\right)\mathbf{W_m^3} + \mathbf{b_m^2}$$

where $H_n^d$ denotes the head number, $r_{ij}$ is the embedding of the relation between nodes $i$ and $j$, $\mathbf{W_m^{1,2,3}}$ and $\mathbf{b_m^{1,2}}$ are weight matrices and biases.

*3) Dual Syntactic Information Fusion:* Now the Position-aware GCN has captured the important local syntactic information and the Relational MHA has captured the important global relational information. To integrate them together and let them compensate for each other, we concatenate the aspect node representations respectively derived by Position-aware GCN and Relational MHA, then we employ a multi-layer perception (MLP), which can automatically abstract the integrated representation [41], [42], to generate the unified node representation sequence, which include the unified aspect representation $\widetilde{\mathbf{R_a}}$.

### D. Knowledge Re-enhancement

After graph modeling, sufficient syntactic information has been integrated into $\widetilde{\mathbf{R_a}}$. On the one hand, some new clues may be captured by DSG-Net and retained in $\widetilde{\mathbf{R_a}}$. Thus $\widetilde{\mathbf{R_a}}$ may further need more aspect knowledge to collaborate with these new clues to support ASC. On the other hand, as the syntax graph may be imperfectly generated by the parser, some wrong connections and relations may be introduced. In this case, re-integrating some knowledge can help alleviate the influence of the imperfect syntax graph. To this end, we design a knowledge integrating gate (KI Gate) to re-enhance $\widetilde{\mathbf{R_a}}$ with further needed knowledge contained in $\mathbf{r_k^T}$. The function of KI gate is given as:

$$\mathbf{R_a} = \widetilde{\mathbf{R_a}} + \mathbf{r_k^T} * \mathbf{W_k^r}[\widetilde{\mathbf{R_a}}, \mathbf{r_k^T}] \qquad (10)$$

where $\mathbf{W_k^r}$ is weight matrix. Here $\widetilde{\mathbf{R_a}}$ and $\mathbf{r_k^T}$ produces a gate scalar rather than a gate vector. There is no subsequent contextualizing module thus $\mathbf{r_k^T}$ can be directly integrated into $\widetilde{\mathbf{R_a}}$ without fine-tuning for adaption. In Sec. IV-F, we investigate the effect of different knowledge gates used here.

### E. Aspect-related Semantics Aggregation

Here we employ an Aspect-to-Context Attention (A2C Att) mechanism to aggregate the aspect-related semantics retained

in all hidden states into a final representation $\mathbf{R_f}$. Similar to A2D Att, A2C Att can be formulated as:

$$\beta_i = \texttt{SoftMax}(\mathcal{F}(h_c^i, \mathbf{R_a})) \tag{11}$$

$$\mathcal{F}(h_c^i, R_a) = (\mathbf{W_{ac}}\, h_c^i + \mathbf{b_{ac}})\, (\mathbf{R_a})^{\mathsf{T}} \tag{12}$$

$$\mathbf{R_f} = \sum_{i=1}^{N} \alpha_i h_c^i \tag{13}$$

where $\mathbf{W_{ac}}$ and $\mathbf{b_{ac}}$ are weight and bias.

### F. Sentiment Classification

We concatenate $\mathbf{R_f}$ with $\mathbf{h_{cls}}$ and then fed the final vector into a linear layer, which is followed by a $\texttt{SoftMax}$ classifier for prediction:

$$\mathbf{P} = \texttt{SoftMax}(\mathbf{W_p}[\mathbf{h_{cls}}, \mathbf{R_f}] + \mathbf{b_p}) \tag{14}$$

where $\mathbf{P}$ is the predicted sentiment distribution, $\mathbf{W_p}$ and $\mathbf{b_p}$ are weight matrix and bias. The cross-entropy loss function is adopted for model training.

There are two reasons why we introduce $\mathbf{h_{cls}}$ here. First, this can add a skip connection to BERT, shortening its loss back-propagation path to facilitate training. The second is for robustness. Possibly the syntax graphs are imperfect and the integrated knowledge contains noise. Hence $\mathbf{h_{cls}}$ serves as a reference and makes the whole model more robust.

## IV. EXPERIMENTS

### A. Datasets

We conduct experiments on three popular datasets for the ASC task: Lap14 and Res14 datasets are from SemEval 2014 task 4 [7], and Res15 dataset is from SemEval 2015 task 12 [43]. The statistics of all datasets are presented in Table I.

TABLE I
DATASET STATISTICS OF THE THREE DATASETS.

| Dataset | Positive | | Neutral | | Negative | |
|---|---|---|---|---|---|---|
| | Train | Test | Train | Test | Train | Test |
| Lap14 | 994 | 341 | 464 | 169 | 870 | 128 |
| Res14 | 2164 | 728 | 637 | 196 | 807 | 196 |
| Res15 | 912 | 326 | 36 | 34 | 256 | 182 |

### B. Experiment Setup

We adopt the BERT-base uncased version [18]. We train our model using Adam optimizer [44] with default configuration. The hyper-parameters are listed in Table II. Accuracy (Acc) and Macro-F1 (F1) are adopted as evaluation metrics. As there is no official validation set, following previous works [12], [13], [45], we run our model three times with random initialization and report the average results on test sets, as shown in Table III. And to compare with the works reporting best results, we also report the best results on test sets, as shown in Table IV. All computations are done on an NVIDIA Quadro RTX 6000 GPU.

TABLE II
SETTING OF HYPER-PARAMETERS.

| Hyper-params | Dataset | | |
|---|---|---|---|
| | Lap14 | Res14 | Res15 |
| learning rate | $1 \times 10^{-5}$ | $5 \times 10^{-5}$ | $3 \times 10^{-5}$ |
| batch size | 32 | 32 | 32 |
| dropout rate | 0.3 | 0.3 | 0.3 |
| $d_e$ | 768 | 768 | 768 |
| $d_s$ | 256 | 256 | 128 |
| $H_n^s$ | 3 | 3 | 6 |
| $H_n^d$ | 2 | 4 | 6 |
| $\mathbf{T}$ | 4 | 4 | 2 |
| GCN layer number | 2 | 2 | 2 |

### C. Compared Baselines

According to what kinds of external information are utilized, we divide the baselines into several group:

1) No external information is used:
   - IAN [9] separately encodes the aspect and context, then model their interactions using an interactive attention mechanism.

2) External corpus is used:
   - PRET+MULT [23] first pre-trains the model on document-level task, then trains the model on both document-level sentiment classification and ASC in the multi-task learning framework.
   - TransCap [32] utilizes a devised aspect-based capsule network to transfer knowledge from document-level task to aspect-level task.

3) Syntax Graph is used:
   - ASGCN [12] employs a GCN to encode the syntax graph for capturing local syntactic information.
   - BiGCN [45] convolutes over hierarchical syntactic and lexical graphs to encode not only original syntactic information but also the corpus level word co-occurrence information.

4) BERT encoder is used:
   - BERT-SPC [18] takes the same input as our model and use $h_{cls}$ for sentiment classification.
   - AEN-BERT [46] adopts BERT encoder and uses the attentional encoder network to model the interactions between the aspect and context.

5) Both of syntax graph and BERT encoder are used:
   - R-GAT+BERT [15] use the relational graph attention network to aggregate the global relational information from all context word into the aspect node representation.
   - DGEDT-BERT [13] employs a dual-transformer network to model the interactions between the flat textual knowledge and dependency graph empowered knowledge.
   - A-KVMN+BERT [47] uses a key-value memory network to leverage not only word-word relations but also their dependency types.
   - BERT+T-GCN [48] leverages the dependency types in T-GCN and use an attentive layer ensemble to learn

TABLE III
PERFORMANCES COMPARISONS OF AVERAGE RESULTS WITH RANDOM INITIALIZATION. $\mathcal{K}, \mathcal{B}, \mathcal{T}$ AND $\mathcal{G}$ DENOTE THE MODEL LEVERAGES ASPECT $\mathcal{K}$NOWLEDGE, $\mathcal{B}$ERT, EXTRA $\mathcal{T}$RAINING CORPUS AND SYNTAX $\mathcal{G}$RAPH, RESPECTIVELY. BEST RESULTS ARE IN **BOLD** AND PREVIOUS SOTA RESULTS ARE UNDERLINED. * DENOTES THAT WE PRODUCE THE RESULTS USING THEIR ORIGINAL SOURCE CODES. $^\dagger$ INDICATES KaGRMN-DSG SIGNIFICANTLY OUTPERFORMS BASELINES UNDER T-TEST ($p < 0.01$).

| External Information | | | Model | Lap14 | | Res14 | | Res15 | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Acc | F1 | Acc | F1 | Acc | F1 |
| — | — | — | IAN [9] | 72.05 | 67.38 | 79.26 | 70.09 | 78.54 | 52.65 |
| — | $\mathcal{T}$ | — | PRET+MULT [23] | 71.15 | 67.46 | 79.11 | 69.73 | 81.30 | 68.74 |
| — | $\mathcal{T}$ | — | TransCap [32] | 73.51 | 69.81 | 79.55 | 71.41 | - | - |
| — | — | $\mathcal{G}$ | ASGCN [12] | 75.55 | 71.05 | 80.77 | 72.02 | 79.89 | 61.89 |
| — | — | $\mathcal{G}$ | BiGCN [45] | 74.59 | 71.84 | 81.97 | 73.48 | 81.16 | 64.79 |
| — | $\mathcal{B}$ | — | BERT-SPC* [18] | 78.47 | 73.67 | 84.94 | 78.00 | 83.40 | 65.00 |
| — | $\mathcal{B}$ | — | AEN-BERT [46] | 79.93 | 76.31 | 83.12 | 73.76 | - | - |
| — | $\mathcal{B}$ | $\mathcal{G}$ | R-GAT+BERT* [15] | 79.31 | 75.40 | 86.10 | <u>80.04</u> | 83.95 | 69.47 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | DGEDT-BERT [13] | 79.8 | 75.6 | <u>86.3</u> | 80.0 | 84.0 | 71.0 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | A-KVMN+BERT* [47] | 79.20 | 75.76 | 85.89 | 78.29 | 83.89 | 67.88 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | BERT+T-GCN* [48] | <u>80.56</u> | <u>76.95</u> | 85.95 | 79.40 | <u>84.81</u> | <u>71.09</u> |
| $\mathcal{K}$ | $\mathcal{B}$ | $\mathcal{G}$ | KaGRMN-DSG (Ours) | **81.87**$^\dagger$ | **78.43**$^\dagger$ | **87.35**$^\dagger$ | **81.21**$^\dagger$ | **86.59**$^\dagger$ | **74.46**$^\dagger$ |
| | | | Our Improvements | **1.62%** | **1.92%** | **1.22%** | **1.46%** | **2.10%** | **4.74%** |

TABLE IV
PERFORMANCES COMPARISONS OF BEST RESULTS. $\mathcal{K}, \mathcal{B}, \mathcal{T}$ AND $\mathcal{G}$ DENOTE THE MODEL LEVERAGES ASPECT $\mathcal{K}$NOWLEDGE, $\mathcal{B}$ERT, EXTRA $\mathcal{T}$RAINING CORPUS AND SYNTAX $\mathcal{G}$RAPH, RESPECTIVELY. BEST RESULTS ARE IN **BOLD** AND PREVIOUS SOTA RESULTS ARE UNDERLINED. * DENOTES THAT WE PRODUCE THE RESULTS USING THEIR ORIGINAL SOURCE CODES.

| External Information | | | Model | Lap14 | | Res14 | | Res15 | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Acc | F1 | Acc | F1 | Acc | F1 |
| — | $\mathcal{B}$ | — | BERT-SPC* [18] | 78.84 | 73.95 | 85.80 | 78.48 | 83.76 | 68.33 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | SAGAT [14] | 80.37 | 76.94 | 85.08 | 77.94 | - | - |
| — | $\mathcal{B}$ | $\mathcal{G}$ | KGCapsAN-BERT [49] | 79.47 | 76.61 | 85.36 | 79.00 | - | - |
| — | $\mathcal{B}$ | $\mathcal{G}$ | R-GAT+BERT* [15] | 79.46 | 75.75 | <u>86.61</u> | <u>80.78</u> | 84.13 | 71.12 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | A-KVMN+BERT [47] | 79.78 | 76.14 | 85.98 | 77.94 | 84.14 | 68.49 |
| — | $\mathcal{B}$ | $\mathcal{G}$ | BERT+T-GCN [48] | <u>80.88</u> | <u>77.03</u> | 86.16 | 79.95 | <u>85.26</u> | <u>71.69</u> |
| $\mathcal{K}$ | $\mathcal{B}$ | $\mathcal{G}$ | KaGRMN-DSG (Ours) | **82.13** | **79.42** | **87.68** | **81.98** | **87.08** | **75.34** |
| | | | Our Improvements | **1.55%** | **3.10%** | **1.24%** | **1.49%** | **2.13%** | **5.09%** |

the comprehensive representation from different T-GCN layers.
- SAGAT [14] utilizes graph attention network and BERT to fully obtain both syntax and semantic information.
- KGCapsAN-BERT [49] utilizes multi-prior knowledge to guide the capsule attention process and use a GCN-based syntactic layer to integrate the syntactic knowledge.

And we label all models with what kinds of external information they leverage, as shown in Table III and Table IV.

### D. Main Results

The performance comparison of all models on average scores is shown in Table III, and the comparison on best scores is shown in Table IV. We can observe that: Syntax graphs, external training corpus, and BERT can all improve ASC. Especially, simple BERT-SPC significantly outperforms all models that do not adopt BERT, even if some of them leverage syntax graph and external training corpus. This shows the power of pre-trained language models on ASC. And combining BERT and syntactic information can further improve results as sufficient semantics captured by BERT and the syntactic information conveyed by syntax graphs can cooperate to assist ASC. However, all baselines do not leverage aspect

knowledge and only consider either local syntactic information or global relational information. As a result, their derived aspect representation lack some important clues of aspect and their captured syntactic information is insufficient, leading to their inferior performance compared to our KaGRMN-DSG model.

We obtain consistent improvements over baselines in terms of Acc and F1 on all datasets, achieving new state-of-the-art results. On average results, our KaGRMN-DSG overpasses previous best results by 1.92%, 1.46%, and 4.74% in terms of Macro-F1 on Lap14, Res14, and Res15 datasets respectively. On best results, KaGRMN-DSG overpasses previous best results by 3.10%, 1.49%, and 5.09% in terms of Macro-F1 on Lap14, Res14, and Res15 datasets respectively. The improvements are contributed by the superiorities of KaGR-MN, which effectively leverage beneficial aspect knowledge, and DSG-Net, which combines GCN and Relational MHA to capture sufficient syntactic information.

### E. Ablation Study

We empirically analyze KaGRMN-DSG and prove the necessity of every component by conducting an ablation study, whose results are shown in Table V. In this section we answer the following research questions (RQs):

TABLE V
RESULTS OF ABLATION STUDY.

| Variants | Lap14 | Res14 | Res15 |
|---|---|---|---|
| | Acc | Acc | Acc |
| $M_0$: KaGRMN-DSG (full model) | **81.87** | **87.35** | **86.59** |
| $M_1$: w/o Aspect Knowledge (descriptions are replaced with aspects) | 80.30 ($\downarrow$ 1.57) | 86.43 ($\downarrow$ 0.92) | 85.24 ($\downarrow$ 1.35) |
| $M_2$: only KaGRMN (w/o DSG-Net + KI Gate + A2C Att) | 80.72 ($\downarrow$ 1.15) | 86.55 ($\downarrow$ 0.80) | 85.36 ($\downarrow$ 1.23) |
| $M_3$: w/o DSG-Net | 80.56 ($\downarrow$ 1.31) | 86.43 ($\downarrow$ 0.92) | 85.56 ($\downarrow$ 1.03) |
| $M_4$: w/o Relational MHA | 80.98 ($\downarrow$ 0.89) | 86.67 ($\downarrow$ 0.68) | 85.79 ($\downarrow$ 0.8) |
| $M_5$: w/o Position-aware GCN | 81.03 ($\downarrow$ 0.84) | 86.76 ($\downarrow$ 0.59) | 85.67 ($\downarrow$ 0.92) |
| $M_6$: w/o KI Gate | 81.09 ($\downarrow$ 0.78) | 87.11 ($\downarrow$ 0.24) | 85.79 ($\downarrow$ 0.80) |
| $M_7$: w/o A2C Att | 81.50 ($\downarrow$ 0.37) | 87.00 ($\downarrow$ 0.35) | 86.41 ($\downarrow$ 0.18) |
| $M_8$: w/o A2D Att | 80.93 ($\downarrow$ 0.94) | 86.70 ($\downarrow$ 0.65) | 85.61 ($\downarrow$ 0.98) |
| $M_9$: w/o Self MHA | 80.36 ($\downarrow$ 1.51) | 86.46 ($\downarrow$ 0.89) | 85.24 ($\downarrow$ 1.35) |

*Effect of Aspect Knowledge.* To study the pure impact of aspect knowledge, we devise two variants: $M_1$ and $M_2$. In $M_1$, the original description is replaced with the aspect itself. In this case, there is no aspect knowledge available for KaGR-MN and its function becomes modeling the interactions between the aspect and context. Surprisingly, even without knowledge, $M_2$ can obtain promising results. We attribute this to the advanced architecture and effective functions of KaGR-MN, in which aspect and context are separately encoded and their interactions are effectively modeled by KaGR-MN. On the other hand, the performance degradation of $M_0$ convincingly demonstrates the pure improvements contributed by the aspect knowledge conveyed by aspect descriptions. In $M_2$, DSG-Net, KI Gate, and A2C Att are all removed, so $M_2$ has a BERT+KaGR-MN architecture and the final aspect representation is used for prediction. $M_2$ consistently outperforms baselines, proving that KaGR-MN can derive a good enough aspect representation in which the clues for aspect sentiment reasoning are retained. Along time steps, recurrently leveraging aspect knowledge, KaGR-MN can capture more and more beneficial clues, semantics and dependencies then retain them in aspect representation and context memories. And effectively utilizing beneficial aspect knowledge is the key advantage of our method compared with previous works.

*Effect of Syntactic Information.* The results gap of $M_3$ and $M_0$ shows the improvement DSG-Net achieves by cooperating with the aspect knowledge. These results validate the advantages of combining both kinds of syntactic information to capture sufficient syntactic information. We then study the effects of Position-aware GCN and Relational MHA. We can observe that both $M_4$ and $M_5$ perform worse than $M_0$, proving both the local syntactic information and global relational information should be captured for ASC. In previous works, only either one of them is considered, leading to insufficient syntactic information. In contrast, our model marries them and lets them compensate for each other, sufficiently capturing syntactic clues.

*Effect of Knowledge Integration Gate.* Without KI Gate, $M_6$ obtains worse results than $M_0$. This indicates that after DSG-Net, some aspect knowledge is further needed and KI Gate is efficient to re-enhance the final aspect representation with the needed knowledge.

*Effect of Aspect-to-Context Attention.* In $M_7$, the final aspect representation is used for prediction. We can find that $M_7$ has limited performance degradation compared to $M_0$. This proves that although previous modules can discover and extract clues for ASC, there are still important clues contained in non-aspect hidden states rather than final aspect representation. Hence it is necessary to employ A2C Att to aggregate the aspect-related semantics in all hidden states into the final representation.

*Effect of A2D Att and Self MHA in KaGR-MN Cell.* The significant performance decrease of $M_8$ shows that A2D Att is indispensable to dynamically summarize the specifically needed aspect knowledge from $\mathbf{M_D}$. Without Self MHA, the integrated knowledge in aspect representation can not be contextualized and context memories cannot be updated. As a result, $M_9$ performs much worse than $M_0$.

### F. Investigation on Knowledge Gates

KaGRMN-DSG has two different knowledge gates (AdaKI and KI) for knowledge integrating. Here we empirically investigate these two knowledge gates by testing their four different settings. The results are shown in Table VI. We can find that $M_{10}$ and $M_{12}$ have slight decreases in performances when respectively compared with $M_0$ and $M_{11}$. This is because KI Gate can preserve the knowledge in $\mathbf{r_k^T}$ while AdaKI Gate may lose some knowledge when adapting to the semantic space of $\widetilde{\mathbf{R_a}}$. $M_{11}$ and $M_{12}$ perform much worse than $M_0$ and $M_{10}$. This is because the semantic space adaption of AdaKI Gate in KaGR-MN can maintain the semantics consistency of $\mathbf{r_a^{t*}}$ and $\mathbf{M_C^{t-1}}$, which is crucial for subsequent knowledge contextualizing.

### G. Impact of Time Step Number $\mathbf{T}$

We plot the performance trends of KaGRMN-DSG with increasing $\mathbf{T}$ on the three datasets, as presented in Fig. 3. We can observe that the performances show a trend of increases at first and then decreases. And the best result is obtained when $\mathbf{T}$ is 2 or 3 for Res15 and 4 for Lap14 and Res14. This shows that appropriately increasing $\mathbf{T}$ can gradually improve the results, which is consistent with our expectation. This can also prove the effectiveness of the recurrent manner of KaGR-MN. However, too large $\mathbf{T}$ leads to inferior performances,

TABLE VI
RESULTS OF DIFFERENT KNOWLEDGE GATE SETTINGS.

| Variants | Gate 1 | Gate 2 | Lap14 Acc | Res14 Acc | Res15 Acc |
|---|---|---|---|---|---|
| $M_0$ | AdaKI | KI | **81.87** | **87.35** | **86.59** |
| $M_{10}$ | AdaKI | AdaKI | 81.50 ($\downarrow$ 0.37) | 87.05 ($\downarrow$ 0.30) | 85.98 ($\downarrow$ 0.61) |
| $M_{11}$ | KI | KI | 81.09 ($\downarrow$ 0.78) | 86.73 ($\downarrow$ 0.62) | 85.36 ($\downarrow$ 1.23) |
| $M_{12}$ | KI | AdaKI | 81.03 ($\downarrow$ 0.84) | 86.58 ($\downarrow$ 0.77) | 85.24 ($\downarrow$ 1.35) |



Fig. 3. Impact of the time step number **T**

TABLE VII
CASES DEMONSTRATION. [**N**, **P**, **O**] DENOTES PREDICTED SENTIMENT DISTRIBUTION: [NEGATIVE, POSITIVE, NEUTRAL].

| Case | [**N**, **P**, **O**] |
|---|---|
| 1. **C**: The **[Mountain Lion OS]**[A] is not hard to figure out if you are familiar with Microsoft Windows. <br> **D**: OS X Mountain Lion is ... Apple Inc.'s desktop and server operating system ... | $M_0$: [0.0, **0.999**$^\checkmark$, 0.001] <br> $M_1$: [0.01, 0.49$^\times$, 0.5] |
| 2. **C**: On start up it asks endless questions just so **[iTune]**[A] can sell you more of their products. <br> **D**: iTunes is a media player, media library, Internet radio broadcaster, mobile device management utility ... | $M_0$: [**0.57**$^\checkmark$, 0.41, 0.02] <br> $M_1$: [0.03, 0.67$^\times$, 0.30] |
| 3. **C**: While the **[smoothies]**[A] are a little big for me, the fresh juices are the best i have ever had! <br> **D**: A smoothie is a drink made from pureed raw fruit and/or vegetables, typically using a blender ... | $M_0$: [**0.62**$^\checkmark$, 0.0, 0.38] <br> $M_1$: [0.02, 0.97$^\times$, 0.01] |
| 4. **C**: All the various Greek and Cypriot dishes are excellent, but the **[gyro]**[A] is the reason to come – if you don't eat one your trip was wasted. **D**: A gyro or gyros is a Greek dish made from meat cooked on a ... | $M_0$: [0.02, **0.98**$^\checkmark$, 0.0] <br> $M_1$: [0.88$^\times$, 0.11, 0.01] |

which is also consistent with our expectation. One possible explanation is that too much knowledge integrated into the aspect representation and context memories will harm their original contextual information. Another is that too many recurrent steps will lead to overfitting on training sets.

### H. Case Study

We show some cases in Table VII. Note that the only difference between KaGRMN-DSG ($M_0$) and $M_1$ is that the input **D** in $M_1$ is replaced with **A**. We can observe that $M_0$ can accurately predict the correct labels in all cases, while $M_1$ fails all cases although its overall performance is promising (as shown in Table V)

Without leveraging aspect knowledge, the aspect representation and semantics derived by $M_1$ are inadequate. As shown in Table VIII, BERT cannot capture the exact meanings and properties of *Mountain Lion OS* and *iTune*, although it is one of the strongest language models. In Case 1, $M_1$ regards *Mountain Lion OS* as 'lion' which is 'dangerous'. Then considering 'not

hard', $M_1$ is confused on P and O. In contrast, leveraging aspect knowledge, $M_0$ captures the exact meaning: an operating system. Then considering the aspect-related semantics ('not hard'), $M_0$ correctly predicts P. In Case 2, the aspect sentiment expression is a little obscure as there are no explicit sentiment trigger words (e.g. delicious, good, expensive). Even if $M_1$ captures aspect-related context semantics, it fails due to the lack of property information of *iTune*. Thanks to the integrated aspect knowledge, $M_0$ is aware that *iTune* is primarily used for media playing rather than selling products, thus correctly predicts N.

Looking into Case 3 and Case 4, we can find that due to the lack of aspect knowledge, $M_1$ is prone to be affected by some misleading sentiment trigger words: 'best' in case 3 and 'but' in case 4. The reason why $M_0$ wins $M_1$ is that $M_0$ can combine the aspect knowledge and the aspect-related semantics together to capture the correct clues for ASC.

TABLE VIII
MISUNDERSTANDING FROM BERT PRESENTED BY SEMANTIC COSINE
SIMILARITY ($S$). $v$ IS THE AVERAGE OF ENTITY'S HIDDEN STATES. $a_i$
DENOTES THE **A** IN CASE $i$.

| Entity ($e$) | $S(v_e, v_{a_1})$ | Entity ($e$) | $S(v_e, v_{a_2})$ |
|---|---|---|---|
| lion | 0.8516 | media player | 0.4720 |
| mountain | 0.7997 | radio broadcaster | 0.5887 |
| operating system | 0.6826 | software | 0.7051 |
| dangerous animal | 0.8272 | utility | 0.6982 |

TABLE IX
COMPARISON OF TRAINING TIME AND INFERENCE TIME (PER SAMPLE) AS
WELL AS THE AVG F1 ON THE THREE DATASETS.

| Models | Training Time↓ | Inference Time↓ | Avg F1↑ |
|---|---|---|---|
| BERT-SPC | **0.007309**s | **0.002219**s | 73.59% |
| BERT+T-GCN | 0.033835s | 0.003350s | 76.22% |
| KaGRMN-DSG | 0.015333s | 0.004208s | **78.91%** |

*I. Computation Time Analysis*

The comparison of time costing and avg F1 of BERT-SPC, BERT+T-GCN and our KaGRMN-DSG model is shown in Table IX. We can find that although our model demands more training time and inference time than BERT-SPC, it overpasses BERT-SPC on avg F1 by a large margin (6.3%). As for BERT+T-GCN, which is the best-performing baseline, although it costs lightly less inference time than our KaGRMN-DSG, it costs much more time for training, and more importantly, its performance is significantly inferior to us. Additionally, since *Local Syntactic Information Modeling* and *Global Relational Information Modeling* both take the output of KaGRMN as input, they can be parallelized theoretically, so the training time and inference time of our KaGRMN-DSG model can be further reduced in practice. In a word, our model may cost more time for training and inference than some baseline models, but it is worthy considering the significant performance improvement.

## V. CONCLUSION

In this paper, we point out the two challenges encountering existing ASC models and we therefore propose a novel KaGRMN-DSG model to end-to-end embed and leverage aspect knowledge, then capture sufficient syntactic information by marrying both kinds of syntactic information. In our model, the integrated beneficial aspect knowledge and sufficient syntactic information can effectively cooperate, yielding new state-of-the-art results.

Future directions include exploring the visual knowledge of aspects, as well as designing deeper and more sufficient dual syntactic interaction to let the two kinds of syntactic information interact with each other in their respective modeling processes.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] K. Schouten and F. Frasincar, "Survey on aspect-level sentiment analysis," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 3, pp. 813–830, March 2016.
[2] L. Zhang, S. Wang, and B. Liu, "Deep learning for sentiment analysis: A survey," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, p. e1253, 2018.
[3] L. Chen, M. Li, W. Su, M. Wu, K. Hirota, and W. Pedrycz, "Adaptive feature selection-based adaboost-knn with direct optimization for dynamic emotion recognition in human–robot interaction," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 5, no. 2, pp. 205–213, 2021.
[4] Y. Xiao, H. Zhao, and T. Li, "Learning class-aligned and generalized domain-invariant representations for speech emotion recognition," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 4, no. 4, pp. 480–489, 2020.
[5] S. Poria, I. Chaturvedi, E. Cambria, and A. Hussain, "Convolutional mkl based multimodal emotion recognition and sentiment analysis." in *ICDM*, F. Bonchi, J. Domingo-Ferrer, R. Baeza-Yates, Z.-H. Zhou, and X. Wu, Eds. IEEE Computer Society, 2016, pp. 439–448. [Online]. Available: http://dblp.uni-trier.de/db/conf/icdm/icdm2016.html#PoriaCCH16
[6] D. Bertero, F. B. Siddique, C.-S. Wu, Y. Wan, R. H. Y. Chan, and P. Fung, "Real-time speech emotion and sentiment recognition for interactive dialogue systems," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics, Nov. 2016, pp. 1042–1047. [Online]. Available: https://aclanthology.org/D16-1110
[7] M. Pontiki, D. Galanis, J. Pavlopoulos, H. Papageorgiou, I. Androutsopoulos, and S. Manandhar, "SemEval-2014 task 4: Aspect based sentiment analysis," in *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*. Dublin, Ireland: Association for Computational Linguistics, 2014, pp. 27–35. [Online]. Available: https://www.aclweb.org/anthology/S14-2004
[8] L. Dong, F. Wei, C. Tan, D. Tang, M. Zhou, and K. Xu, "Adaptive recursive neural network for target-dependent Twitter sentiment classification," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Baltimore, Maryland: Association for Computational Linguistics, 2014, pp. 49–54. [Online]. Available: https://www.aclweb.org/anthology/P14-2009
[9] D. Ma, S. Li, X. Zhang, and H. Wang, "Interactive attention networks for aspect-level sentiment classification," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, C. Sierra, Ed. ijcai.org, 2017, pp. 4068–4074. [Online]. Available: https://doi.org/10.24963/ijcai.2017/568
[10] F. Fan, Y. Feng, and D. Zhao, "Multi-grained attention network for aspect-level sentiment classification," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, 2018, pp. 3433–3442. [Online]. Available: https://www.aclweb.org/anthology/D18-1380
[11] J. Tang, Z. Lu, J. Su, Y. Ge, L. Song, L. Sun, and J. Luo, "Progressive self-supervised attention learning for aspect-level sentiment analysis," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, 2019, pp. 557–566. [Online]. Available: https://www.aclweb.org/anthology/P19-1053
[12] C. Zhang, Q. Li, and D. Song, "Aspect-based sentiment classification with aspect-specific graph convolutional networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, 2019, pp. 4568–4578. [Online]. Available: https://www.aclweb.org/anthology/D19-1464
[13] H. Tang, D. Ji, C. Li, and Q. Zhou, "Dependency graph enhanced dual-transformer structure for aspect-based sentiment classification," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, 2020, pp. 6578–6588. [Online]. Available: https://www.aclweb.org/anthology/2020.acl-main.588
[14] L. Huang, X. Sun, S. Li, L. Zhang, and H. Wang, "Syntax-aware graph attention network for aspect-level sentiment classification," in *Proceedings of the 28th International Conference on Computational Linguistics*. Barcelona, Spain (Online): International Committee on Computational Linguistics, 2020, pp. 799–810. [Online]. Available: https://www.aclweb.org/anthology/2020.coling-main.69

[15] K. Wang, W. Shen, Y. Yang, X. Quan, and R. Wang, "Relational graph attention network for aspect-based sentiment analysis," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, 2020, pp. 3229–3238. [Online]. Available: https://www.aclweb.org/anthology/2020.acl-main.295

[16] K. Sun, R. Zhang, S. Mensah, Y. Mao, and X. Liu, "Aspect-level sentiment analysis via convolution over dependency tree," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, 2019, pp. 5679–5688. [Online]. Available: https://www.aclweb.org/anthology/D19-1569

[17] P. Chen, Z. Sun, L. Bing, and W. Yang, "Recurrent attention network on memory for aspect sentiment analysis," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, 2017, pp. 452–461. [Online]. Available: https://www.aclweb.org/anthology/D17-1047

[18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 4171–4186. [Online]. Available: https://www.aclweb.org/anthology/N19-1423

[19] L. Jiang, M. Yu, M. Zhou, X. Liu, and T. Zhao, "Target-dependent Twitter sentiment classification," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics, 2011, pp. 151–160. [Online]. Available: https://www.aclweb.org/anthology/P11-1016

[20] S. Mohammad, S. Kiritchenko, and X. Zhu, "NRC-Canada: Building the state-of-the-art in sentiment analysis of tweets," in *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*. Atlanta, Georgia, USA: Association for Computational Linguistics, 2013, pp. 321–327. [Online]. Available: https://www.aclweb.org/anthology/S13-2053

[21] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM for aspect-level sentiment classification," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics, 2016, pp. 606–615. [Online]. Available: https://www.aclweb.org/anthology/D16-1058

[22] D. Tang, B. Qin, X. Feng, and T. Liu, "Effective LSTMs for target-dependent sentiment classification," in *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. Osaka, Japan: The COLING 2016 Organizing Committee, 2016, pp. 3298–3307. [Online]. Available: https://www.aclweb.org/anthology/C16-1311

[23] R. He, W. S. Lee, H. T. Ng, and D. Dahlmeier, "Exploiting document knowledge for aspect-level sentiment classification," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 579–585. [Online]. Available: https://www.aclweb.org/anthology/P18-2092

[24] B. Huang and K. Carley, "Parameterized convolutional neural networks for aspect level sentiment classification," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, 2018, pp. 1091–1096. [Online]. Available: https://www.aclweb.org/anthology/D18-1136

[25] W. Xue and T. Li, "Aspect based sentiment analysis with gated convolutional networks," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 2514–2523. [Online]. Available: https://www.aclweb.org/anthology/P18-1234

[26] S. Wang, S. Mazumder, B. Liu, M. Zhou, and Y. Chang, "Target-sensitive memory networks for aspect sentiment classification," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 957–967. [Online]. Available: https://www.aclweb.org/anthology/P18-1088

[27] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. [Online]. Available: https://openreview.net/forum?id=rJXMpikCZ

[28] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: https://openreview.net/forum?id=SJU4ayYgl

[29] M. S. Schlichtkrull, T. N. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling, "Modeling relational data with graph convolutional networks." in *ESWC*, ser. Lecture Notes in Computer Science, A. Gangemi, R. Navigli, M.-E. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai, and M. Alam, Eds., vol. 10843. Springer, 2018, pp. 593–607. [Online]. Available: http://dblp.uni-trier.de/db/conf/esws/eswc2018.html#SchlichtkrullKB18

[30] R. Wang, W. Ji, and B. Song, "Durable relationship prediction and description using a large dynamic graph." *World Wide Web*, vol. 21, no. 6, pp. 1575–1600, 2018. [Online]. Available: http://dblp.uni-trier.de/db/journals/www/www21.html#WangJS18

[31] B. Huang and K. Carley, "Syntax-aware aspect level sentiment classification with graph attention networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, 2019, pp. 5469–5477. [Online]. Available: https://www.aclweb.org/anthology/D19-1549

[32] Z. Chen and T. Qian, "Transfer capsule network for aspect level sentiment classification," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, 2019, pp. 547–556. [Online]. Available: https://www.aclweb.org/anthology/P19-1052

[33] B. Xing, L. Liao, D. Song, J. Wang, F. Zhang, Z. Wang, and H. Huang, "Earlier attention? aspect-aware LSTM for aspect-based sentiment analysis," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, S. Kraus, Ed. ijcai.org, 2019, pp. 5313–5319. [Online]. Available: https://doi.org/10.24963/ijcai.2019/738

[34] J. Pennington, R. Socher, and C. Manning, "GloVe: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, 2014, pp. 1532–1543. [Online]. Available: https://www.aclweb.org/anthology/D14-1162

[35] J. Weston, S. Chopra, and A. Bordes, "Memory networks," 2014, cite arxiv:1410.3916. [Online]. Available: http://arxiv.org/abs/1410.3916

[36] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[37] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder–decoder for statistical machine translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, 2014, pp. 1724–1734. [Online]. Available: https://www.aclweb.org/anthology/D14-1179

[38] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 5998–6008. [Online]. Available: https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[39] M. Ravanelli, P. Brakel, M. Omologo, and Y. Bengio, "Light gated recurrent units for speech recognition," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 2, no. 2, pp. 92–102, 2018.

[40] Z. Lin, M. Feng, C. N. dos Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio, "A structured self-attentive sentence embedding," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: https://openreview.net/forum?id=BJC_jUqxe

[41] D.-K. Nguyen and T. Okatani, "Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering." in *CVPR*. IEEE Computer Society, 2018, pp. 6087–6096.

[42] L. Qin, W. Che, Y. Li, M. Ni, and T. Liu, "Dcr-net: A deep co-interactive relation network for joint dialog act recognition and sentiment classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI Press, 2020, pp. 8665–8672. [Online]. Available: https://aaai.org/ojs/index.php/AAAI/article/view/6391

[43] M. Pontiki, D. Galanis, H. Papageorgiou, S. Manandhar, and I. Androut-sopoulos, "SemEval-2015 task 12: Aspect based sentiment analysis," in *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*. Denver, Colorado: Association for Computational Linguistics, Jun. 2015, pp. 486–495.

[44] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: http://arxiv.org/abs/1412.6980

[45] M. Zhang and T. Qian, "Convolution over hierarchical syntactic and lexical graphs for aspect level sentiment analysis," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics, 2020, pp. 3540–3549. [Online]. Available: https://www.aclweb.org/anthology/2020.emnlp-main.286

[46] Y. Song, J. Wang, T. Jiang, Z. Liu, and Y. Rao, "Attentional encoder network for targeted sentiment classification," 2019.

[47] Y. Tian, G. Chen, and Y. Song, "Enhancing aspect-level sentiment analysis with word dependencies," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Online: Association for Computational Linguistics, Apr. 2021, pp. 3726–3739. [Online]. Available: https://www.aclweb.org/anthology/2021.eacl-main.326

[48] ——, "Aspect-based sentiment analysis with type-aware graph convolutional networks and layer ensemble," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, 2021, pp. 2910–2922.

[49] B. Zhang, X. Li, X. Xu, K.-C. Leung, Z. Chen, and Y. Ye, "Knowledge guided capsule attention network for aspect-based sentiment analysis." *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 28, pp. 2538–2551, 2020. [Online]. Available: http://dblp.uni-trier.de/db/journals/taslp/taslp28.html#ZhangLXLCY20

**Ivor W. Tsang** is an IEEE Fellow and the Director of A*STAR Centre for Frontier AI Research (CFAR). Previously, he was a Professor of Artificial Intelligence, at University of Technology Sydney (UTS), and Research Director of the Australian Artificial Intelligence Institute (AAII). His research focuses on transfer learning, deep generative models, learning with weakly supervision, big data analytics for data with extremely high dimensions in features, samples and labels. His work is recognised internationally for its outstanding contributions to those fields. In 2013, Prof Tsang received his ARC Future Fellowship for his outstanding research on big data analytics and large-scale machine learning. In 2019, his JMLR paper "Towards ultrahigh dimensional feature selection for big data" received the International Consortium of Chinese Mathematicians Best Paper Award. In 2020, he was recognized as the AI 2000 AAAI/IJCAI Most Influential Scholar in Australia for his outstanding contributions to the field, between 2009 and 2019. His research on transfer learning was awarded the Best Student Paper Award at CVPR 2010 and the 2014 IEEE TMM Prize Paper Award. In addition, he received the IEEE TNN Outstanding 2004 Paper Award in 2007 for his innovative work on solving the inverse problem of non-linear representations. Recently, Prof Tsang was conferred the IEEE Fellow for his outstanding contributions to large-scale machine learning and transfer learning. Prof Tsang serves as the Editorial Board for the JMLR, MLJ, JAIR, IEEE TPAMI, IEEE TAI, IEEE TBD, and IEEE TETCI. He serves as a Senior Area Chair/Area Chair for NeurIPS, ICML, AAAI and IJCAI, and the steering committee of ACML.



**Bowen Xing** received his B.E. degree and Master degree from Beijing Institute of Technology, Beijing, China, in 2017 and 2020, respectively. He is currently a second year Ph.D student at Australian AI Institute, University of Technology Sydney (UTS). His research focuses on graph neural network, multi-task learning, sentiment analysis and dialog system.