# AR-BERT: Aspect-relation enhanced Aspect-level Sentiment Classification with Multi-modal Explanations

Sk Mainul Islam
IIT Kharagpur, India

Sourangshu Bhattacharya
IIT Kharagpur, India

## ABSTRACT

Aspect level sentiment classification (ALSC) is a difficult problem with state-of-the-art models showing less than 80% macro-F1 score on benchmark datasets. Existing models do not incorporate information on aspect-aspect relations in knowledge graphs (KGs), e.g. DBpedia. Two main challenges stem from inaccurate disambiguation of aspects to KG entities, and the inability to learn aspect representations from the large KGs in joint training with ALSC models. We propose AR-BERT, a novel two-level global-local entity embedding scheme that allows efficient joint training of KG-based aspect embeddings and ALSC models. A novel incorrect disambiguation detection technique addresses the problem of inaccuracy in aspect disambiguation. We also introduce the problem of determining mode significance in multi-modal explanation generation, and propose a two step solution. The proposed methods show a consistent improvement of $2.5 - 4.1$ percentage points, over the recent BERT-based baselines on benchmark datasets. The code is available at `https://github.com/mainuliitkgp/AR-BERT.git`.

## CCS CONCEPTS

• **Information systems** → **Sentiment analysis**; • **Computing methodologies** → **Natural language processing**.

## KEYWORDS

Sentiment Analysis, Knowledge Graph Embedding, Explainable Deep Learning

## 1 INTRODUCTION

Aspect level sentiment classification (ALSC) is an important NLP task [7, 11, 24], where we predict the sentiment portrayed in a *sentence* (also called context) towards an identified *aspect* phrase. Recently, models capturing aspect-specific features, e.g., Transformation Network (TNet) [15], which constructs aspect-specific embedding of context words, or BERT-based models [6], which capture aspect-specific representations of sentences, have outperformed previous sequential prediction models. Other recent improvements

include domain adaptation of the BERT model [27] and incorporating entity relationships within the same sentence using graph convolutional networks [45]. However, existing ALSC methods do not explicitly utilize the relations between aspects, which could potentially lead to better performance.

We observe that many of the aspect phrases, e.g., `Windows 8`, `Mozzarella`, `Taylor Swift`, etc., are mentions of named entities appearing in knowledge graphs (KG), e.g., DBpedia, which encode various entity-entity relations. While some of the aspects may be unseen in the training data, their neighbors (related aspects) in the KG may be abundant. This information can be used to infer important signals from the context sentence, which in turn can help in the prediction of correct polarity. For example, in the sentence, `My laptop with Windows 7 crashed and I did not want Windows 8.`, the aspect `Windows 7` has only 17 examples in the training data. The current state-of-the-art ALSC model [27] wrongly predicts the aspect sentiment as positive. However, its related aspects (1-hop neighbors in the DBpedia KG) have 209 training examples, which can lead to the correct prediction of sentiment. Our primary objective is to incorporate the KG relations between aspects into ALSC models.

The main challenges in model design are: (1) end-to-end training of ALSC models with KG embeddings is infeasible due to large scale of KGs, and (2) most scalable off-the-shelve named entity disambiguation techniques, e.g., wikifier [2] are highly inaccurate. While state of the art named entity disambiguation methods [3, 14] are accurate, they still do not scale to the entire DBpedia KG. We solve the problem of learning aspect representations from large KGs using a two-level graph embedding technique: one corresponding to a higher level *cluster graph*, and another for subgraphs. These embeddings can be efficiently trained jointly along with ALSC models. The problem of inaccurate wikification [2] method for aspect disambiguation, is ameliorated by a novel probing function based detection of incorrect aspect disambiguations. Figure 1 shows the overall architecture of the proposed technique.

The deep multi-modal ALSC model proposed here utilizes information from both text and graph data in an opaque manner, thus reducing the trust in the model predictions. Hence, we seek to design a postthoc global explanation model for predicting multi-modal explanations from both context words, and KG-entities. A key challenge in the prediction of multi-modal explanations is the prediction of importance of the mode for which explanations are being generated, since there may not be any valid explanations from a given mode. We design an explanation model which also predicts the mode importance, and a two-step prodecure for jointly learning the unimodal explanation models, as well as the mode importance predictor. To the best of our knowledge, ours is the first model for jointly predicting multi-modal explanations from text and graph data.
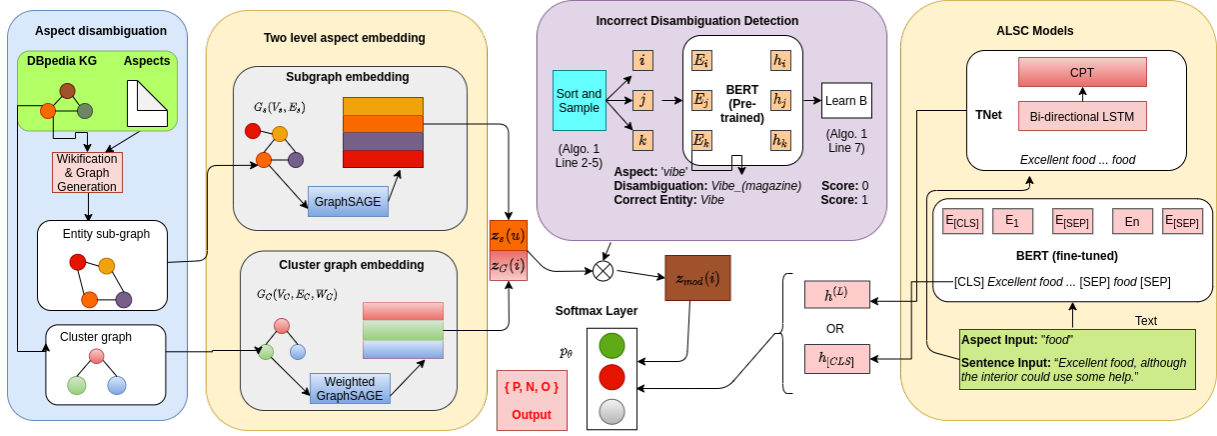
**Figure 1: Architecture of our framework for end-to-end training of ALSC model while incorporating aspect relations from large KGs. Yellow backgound denotes modules trained end-to-end, green backgound denotes input.**

Experimental results show that proposed ALSC models improve the macro-F1 score and accuracy of state-of-the-art ALSC methods on three benchmark datasets by between 2.5% − 4.1%. We also demonstrate that the scarcity of training examples is indeed a factor for the inaccuracy of existing models. We also show that classification accuracy of wrongly disambiguated aspects improves significantly with the disambiguation correction method. Finally, experimental results using the explanation prediction model show that both predicted explanations and significant mode are effective and intuitive. To summarize, our main contributions are: (1) AR-BERT - a scalable aspect-relation enhanced ALSC algorithm. (2) A novel technique for detecting incorrect entity disambiguations. (3) A multi-modal explanation generation model with significant mode detection.

## 2 RELATED WORK

**ALSC with Graph Embedding:** [20] uses memory networks generate aspect representations influenced by other aspects in the same sentence. [42] uses aspect-specific GCN, and [16] uses an "interactive" dependency graph to capture the relations between aspect in a sentence. [32] also encodes information in dependency graph using a transformer-like network. However, none of the above methods can be used at a scale where we can apply it to a knowledge graph like DBpedia. [37],[12] focus on determining aspect specific opinion spans. In addition to the models described in section 3.1, neural network models such as Memory Networks [4, 31, 34], LSTM-based models [19, 35, 43], and Capsule Networks [8] have also been explored for ALSC.

**Knowledge graphs in BERT representations:** Reif et al. [25] measures the word sense similarities using a semantic probe on word embeddings. Hewitt and Manning [10] shows that contextual word embedding incorporates syntactic informations. Broscheit [3] investigates entity knowledge in BERT embedding. [22, 23, 44] propose a promising line of schemes for incorporating entity knowledge in KGs into BERT embeddings. However end-to-end training with these methods has to take entire KG into account, and is expected to be computationally expensive.

[30] modifies the BERT encoder and training procedure in order to incorporate graphs constructed from KG and unstructured text. However, this is not scalable. [17] augments the unstructured text with triples from KG, and trains BERT on the resulting corpora. This technique is also not scalable when trying to capture relations between all entities.

**Explainable models:** For explaining textual and graph information in a unimodal manner popular methods e.g. LIME [26] (extracts important words for a particular prediction) and GNNExplainer [41] (extracts an important subgraph for a particular node classification prediction), cannot be used since they only learn local explanation models. We use the recently proposed methods of concept learning [40] for explanations using textual data, and PGExplainer [18] for explanations using graphs data. While multi-modal explanation generation has been studied in for image and text modes [13, 21], to the best of our knowledge, multi-modal explanation and mode significance has not been studied in the context of textual and graph data.

## 3 EXPLAINABLE ALSC WITH ASPECT RELATION

In this section, we describe our approach for improving the performance of aspect level sentiment classification (ALSC) methods using semantic relations between aspects which can be extracted from Knowledge Graphs (KG), e.g. DBpedia. The key motivation behind our work is that certain aspects are not well represented by examples in the training set, but they have neighboring entities in the KG which have more examples. Hence, the semantic information learned from the neighboring aspect may be transferred to the current aspect through aspect embeddings. For example, in the sentence [However, I can refute that OSX is "FAST"], the aspect OSX has the corresponding DBpedia entity *MacOS*, which has only 7 examples in the training set. However, *MacOS* has a related entity *Microsoft_Windows* which has 37 examples. This leads to the existing BERT-based ALSC method [27] misclassifying this example as positive sentiment polarity based on the context word

*FAST*, whereas our method focusses on the context word *refute* and classifies the example correctly as negative sentiment polarity.

Our method has 2 broad components: (1) disambiguating mentions of aspects (e.g. *OSX*) to entities from a KG (e.g. DBpedia entity *MacOS*), and representing them as an embedding vector, and (2) incorporation of the vector representation of the aspects into state of the art ALSC models, e.g. TNET and BERT, using end-to-end training. Figure 1 describes the overall architecture of our technique. Section 3.1 provides the background on ALSC, sections 3.2, 3.3, and 3.4 describe the proposed ALSC model, AR-BERT, and an incorrect disambiguation detection technique. Section 3.5 describes a novel model for multi-modal explanation generation for AR-BERT.

## 3.1 Background in ALSC

The task of aspect level sentiment classification (ALSC) is to determine the sentiment polarity $y \in \{P, N, O\}$ of an input sentence $w$ for an aspect phrase $w^t$, which is a part of the input sentence. Here, $P$, $N$, and $O$ correspond to positive negative and neutral sentiment respectively. ALSC models take representations of the context $w$, $\vec{x} = (x_1, ..., x_n)$, and that of the aspect $w^t$, $\vec{x}^t = (x_1^t, ..., x_m^t)$ as inputs. Most state-of-the-art ALSC models, including TNET [15], and BERT [6] transform the context representation using an aspect representation to finally arrive at an aspect-specific representation for context words. Here, $n$ denotes the length of the (context) sentence and $m$ denotes the length of the aspect. We briefly describe the architectures of these methods.

TNet [15] consists of three sequential modules (sets of layers): The first module is a Bi-LSTM layer which takes context embeddings $\vec{x}$ and aspect (target) word embeddings $\vec{x}^t$ corresponding to each eaxmple and outputs the contextualized representations $\vec{h}^{(0)} = (h_1^{(0)}(\vec{x}), ..., h_n^{(0)}(\vec{x}))$ and $\vec{h}^t = (h_1^t(\vec{x}^t), ..., h_m^t(\vec{x}^t))$ respectively where $h_i^{(0)}(x), h_j^t(x^t) \in \mathbb{R}^{2D_h}, i \in \{1, ..., n\}, j \in \{1, ..., m\}$. The second module contains $L$ layers of context preserving transformations (CPT) blocks. In each layer $l$ the aspect representation is first transformed into aspect specific representation as $r_i^t = \sum_{j=1}^m h_j^t * \text{SoftMax}(h_i^{(l)}, h_j^t)$, then incorporated into context representation as $\tilde{h_i}^{(l)} = \text{FeedForward}(h_i^{(l)} r_i^t)$, and finally passed into context preserving **LF/AS** block to get the output of next layer: $h_i^{(l+1)} = LF/AS(h_i^{(l)}, \tilde{h_i}^{(l)})$. The third module uses convolution and pooling layers on position-aware encodings to produce a fixed dimensional vector $z$.

BERT has been applied to ALSC by Rietzler et al. [27] and Sun et al. [29], where they model the sentiment classification task as a sequence-pair classification task. The input sentence $(\vec{x})$ and aspect phrase $(\vec{x}^t)$ are encoded as [CLS] $\vec{x}$ [SEP] $\vec{x}^t$ [SEP]. The last layer hidden representation of CLS token $h_{[CLS]} \in \mathbb{R}^{768}$, which is the aspect-aware representation of the input sequence, is used for the downstream classification task. The sentiment polarity distribution is predicted using a feedforward layer with softmax activation and trained using the cross-entropy loss. Recently, SDGCN-BERT [45] has been proposed to capture sentiment dependencies between multiple aspects in a sentence using a graph convolution network. BERT-ADA [27] uses BERT domain-specific language model fine-tuning for ALSC and results in best accuracy on some benchmark datasets. In this paper we build on TNet, BERT, SDGCN-BERT, and

BERT-ADA, to incorporate knowledge from KG. Next, we describe our framework for the scalable incorporation of KG information in ALSC.

## 3.2 Aspect Relation Incorporation from KG

Incorporating aspect relation from KG into ALSC models has two substeps: (1) Aspect to entity mapping and (2) Computation of entity embedding. The first step involves the identification of Wikipedia entities corresponding to an aspect word in a context. This problem is solved by named entity disambiguation (NED) or wikification. We use *wikifier* API [2] for this purpose. Note that, here we use a freely available and computationally efficient method for entity linking, at the cost of accuracy. We partially make up for the loss of accuracy in the posthoc disambiguation correction described in section 3.4.

For learning the entity embeddings (step 2), we use the popular GraphSAGE algorithm [9], which is applicable for both supervised and unsupervised tasks. The entity relation graph is generated using the DBpedia [1] page links knowledge graph, where each vertex is an entity in the DBpedia KG, and an edge is a tuple of the form $< Sub, Pred, Obj >$ where $Sub$ and $Obj$ are the subject and object entities, and $Pred$ is the predicate relation between $Sub$ and $Obj$. However, the whole DBpedia knowledge graph (KG) is too large (with $\sim$ 22 million nodes and $\sim$ 173 million edges) to embed using deep NRL techniques. Another alternative is to consider the subgraph $G$ induced by entities present in the ALSC training dataset only. The problem with this subgraph is that it is disconnected. Hence, the similarity preserving embeddings of entities are only consistent within the connected components of $G$. This may lead to two very different entities $u$ and $v$ accidentally ending up close to each other. In this section, we describe a two-level scalable network embedding technique that scales to DBpedia while avoiding the above-mentioned problems.

*3.2.1 Two-level Aspect Entity Embedding.* The key idea behind two-level aspect embedding (representations) is two use two smaller graphs constructed from the large KG: (1) a cluster graph $G_C(V_C, E_C, W_C)$: which captures the global connectivity structure between clusters of entities, and (2) the subgraph $G_s(V_s, E_s)$ induced by aspects (entities) in the training dataset. Note that since the subgraph $G_s$ can be disconnected, we need a combination of *cluster graph embedding* $\mathbf{z}_C(u)$, and *subgraph embedding*, $\mathbf{z}_s(u)$ for capturing the relations between aspect entity $u$.

**Cluster graph embedding**: The weighted cluster graph $G_C(V_C, E_C, W_C)$ is a compact representation of the KG where each vertex $v \in V_C$ is a cluster of vertices (entities) of the knowledge graph $G = (V, E)$. We use the Louvain hierarchical graph clustering [1] algorithm for clustering the entire knowledge graph. Edge set $E_C$ is calculated as: $(i, j) \in E_C, \forall i, j \in V_C$ if there is a connected pair of KG entities from clusters $i$ and $j$. The weight between clusters $i$ and $j$, $W_C(i, j)$, is calculated as the fraction of actual edges between clusters $i$ and $j$ and the maximum edges possible between the two clusters, i.e, $|i| * |j|$, where $|i|$ is the number of nodes present in cluster $i$. We use a modified GraphSAGE embedding technique to construct the cluster embeddings $\mathbf{z}_C(i), i \in V_C$ of a weighted graph by optimizing

---

[1] https://wiki.dbpedia.org/downloads-2016-10

the following graph based loss function:

$$J_C(z_C(i)) = -log(\sigma(z_C(i)^T z_C(j))) - Q \cdot \mathbb{E}_{k \sim P_n(j)} log(\sigma(-z_C(i)^T z_C(k)))$$

where $z_C(i)$ is the output representation of $i \in V_C$, $\sigma$ is the sigmoid function, $j \in V_C$ is a cluster co-occurring with $i$ on a fixed weighted random walk defined by $W_C(i, j)$, $P_n$ is the negative sampling distribution, $Q$ is the number of negative samples, $k \in V_C$ is a negative sample.

**Subgraph embedding**: The vertex set $V_s$ of the entity-relation subgraph $G_s(V_s, E_s)$ consists of all aspect entities extracted from instances in the training dataset, while the edge set $E_s$ is the subset of induced edges from the original KG. We use the standard GraphSAGE embedding and loss function to construct the *subgraph similarity embedding*, $z_s(u)$ for aspect entity $u$. To preserve the local neighbourhood information as well as global graph structure in the knowledge graph, we use the concatenation of subgraph and cluster graph embedings as our two level entity embedding: $z(u) = [z_C(i); z_s(u)]$, where $u \in V_s$ and $i \in V_C$ such that $u$ is an entitiy in cluster $i$. Figure 1 shows the methods for aspect disambiguation and two-level entity embedding on the left side in the overall scheme.

### 3.3 ALSC with entity relation learning

In this section, we incorporate the concatenated entity embedding proposed above into two state-of-the-art ALSC models: TNet and BERT (described in section 3.1). We propose two ways of incorporating the information contained in entity relations from KG into ALSC: (1) using static embeddings, and (2) by performing end-to-end learning. For incorporation of static embedding in TNet, the final entity embedding $z(u)$ for entity $u$ is concatenated with final layer CPT block output $h^{(L)}$ as $h_{concat}^{(L)} = [h^{(L)}; z(u)]$ and this new aspect specific contextual representation $h_{concat}^{(L)}$ is sent as input to the convolution layer module as described in section 3.1. The final layers and loss function is same as TNet. We call this model **Aspect Relation-TNet (AR-TNet)**. We incorporate the entity embedding $z(u)$ into BERT by concatenating it with representation of CLS token $h_{[CLS]}$, as: $h_{[CLS]_{concat}} = [h_{[CLS]}; z(u)]$. Here $h_{[CLS]}$ is the final aspect-specific sentence representation for an ALSC instance taken from domain specific BERT model (BERT-ADA) [27], and further fine-tuned on ALSC task. We call this model **Aspect Relation-BERT (AR-BERT)**. We also incorporate the static entity embedding $z(u)$ into SDGCN-BERT [45], in an analogous way to train the **Aspect Relation-BERT-S (AR-BERT-S)** model, through finetuning on ALSC. These models are referred to in Table 2 with **'wo end-to-end'** in parenthesis because these models are trained without an end-to-end strategy. The end-to-end training of our proposed models **AR-BERT**, **AR-BERT-S**, and **AR-TNet** (for the base models BERT-ADA, SDGCN-BERT, and TNet, respectively) are discussed in the following section.

**End-to-end learning**: Incorporation of GraphSAGE embeddings into ALSC models provide minor improvements to polarity prediction, since the aspect embeddings are not fine-tuned for the ALSC task. This is achieved with end-to-end training of the aspect embedding and ALSC models. The architecture of our end-to-end models are same as the models proposed with static embeddings above. Hence, for BERT based models, we calculate the final embeddings

for a sentence and aspect pair as: $h_{[CLS]_{concat}} = [h_{[CLS]}; z(u)]$, where $z(u) = [z_C(i); z_s(u)]$. For TNet-based models, $h_{concat}^{(L)} = [h^{(L)}; z(u)]$. For both models, let $\mathcal{L}_{ALSC}$ denote the loss incurred from ALSC training, and $\mathcal{L}_{GS}$ be the loss incurred from GraphSage using the subgraph $G_s = (V_s, E_s)$. We optimize the following loss for joint training:

$$\mathcal{L}_{joint}(\Theta_{ALSC}, \{z_s(u)\}) = \alpha_1 \mathcal{L}_{ALSC} + \alpha_2 \mathcal{L}_{GS}$$

where, $\Theta_{ALSC}$ are all the parameters from ALSC model, and $\{z_s(u)\}$ are subgraph embeddings from $G_s$. We minimize the above loss w.r.t. $\Theta_{ALSC}, \{z_s(u)\}$, while keeping $z_C(i)$ fixed to pre-learned Graph-SAGE embeddings.

### 3.4 Incorrect Disambiguation Detection

Many of the misclassifications using models like BERT-GS, are due to incorrect disambiguation of aspect entities (see section 4.3). In this section, we develop a scalable algorithm for identifying incorrect aspect disambiguations and mitigating their effect by setting the corresponding (modified) embedding to zero vector. We rely on the BERT aspect embedding vectors $h_{[CLS]}$ (called $h$ in this section for brevity) for the same. However, BERT embeddings encode many modalities of information including syntactic dependencies [10], semantic similarities, and entity relations [25]. We propose to use a learned similarity function $S_B(h_i, h_j)$ which captures the entity similarity between two BERT embeddings $h_i$ and $h_j$ of two entity mentions. Hence, following [25], we propose to use the following form of similarity function:

$$S_B(h_i, h_j) = \sigma((B \cdot h_i)^T (B \cdot h_j))$$

where, $B \in R^{dim_B * dim_h}$ is a learned parameter. The parameter $B$ can be thought of as a "probing function" [25], projecting BERT embedding $h$ into a space which only distills out the entity relations.

Algorithm 1, describes the steps for learning the probing function parameter $B$, which extracts entity relational similarities from BERT embeddings, and calculation of the modified embeddings. The key idea is: *aspects which are close in graph embedding space should also have high similarity of BERT embeddings*. The algorithm proceeds by constructing triplets $(i, j, k)$ of aspects where aspects $i$ and $j$ are closer, but $i$ and $k$ are not closer. It then learns $B$ by minimizing the loss: $\sum_{(i,j,k) \in \tau} (S_B(h_i, h_k) - S_B(h_i, h_j))$. For each aspect $i$, and for all it's top $n$ close aspects $j$ and rest far away aspects $k$, we modify it's corresponding concatenated entity embedding as follows:

$$z_{mod}(i) = \begin{cases} \{\vec{0}\}^{dim_h}, & \text{if } S_B(h_i, h_j) - S_B(h_i, h_k) \geq 0 \\ z(i), & \text{otherwise} \end{cases}$$

(1)

$\{\vec{0}\}^{dim_h}$ is the zero vector of dimension $dim_h$. We call ALSC models jointly (end-to-end) trained with these corrected embeddings as: **AR-BERT-idd**, **AR-BERT-S-idd**, and **AR-TNet-idd**; corresponding to base models BERT-ADA, SDGCN-BERT, and TNet.

### 3.5 Multi-modal Explanation Generation

The ALSC models proposed above incorporate semantic and syntactic information from the text data, and aspect relations from the knowledge graphs. Since, the information from these two modes

---

**Algorithm 1** Incorrect Disambiguation Detection

---

**Require:** Set of aspects $\mathcal{A}(\mathcal{D})$ in Dataset $\mathcal{D}$, Graph aspect embeddings $\{z(u)\}, u \in \mathcal{A}(\mathcal{D})$, BERT aspect embeddings $\{\vec{h}(u)\}, u \in \mathcal{A}(\mathcal{D})$

**Ensure:** Probing distance fn. parameter: $B$, Corrected embedding: $z_{mod}(u)$

---
 1: randomly initialize $B$
 2: LIST $\tau \leftarrow \phi$
 3: **for all** $i \in \mathcal{A}(\mathcal{D})$ **do**
 4:    $\mathcal{A}_i(\mathcal{D}) \leftarrow$ sort other aspects $z(u), u \in \mathcal{A}(\mathcal{D})$ in decreasing order of distance to $z(i)$
 5:    Sample $j$ from top $n$ in $\mathcal{A}_i(\mathcal{D})$, and $k$ from rest; emit $\tau \leftarrow \tau \cup (i, j, k)$
 6: **end for**
 7: learn $B$: $\min_B \sum_{(i,j,k) \in \tau} (\mathcal{S}_B(h_i, h_k) - \mathcal{S}_B(h_i, h_j)) + \lambda \|B\|^2$
 8: Calculate $z_{mod}(u)$ using eqn. 1

---

are combined using a deep neural network, for a given test example, the information content in each of the modes is not obvious. In this section, we describe a global posthoc explanation model for generating *multi-modal explanations* for predictions provided by the proposed model architecture. For simplicity, we summarize the proposed architecture into 3 components: (1) the text feature extractor model for input context and aspect $\vec{x}$, $M_t(\vec{x}) : X \to h(\vec{x})$, (2) the graph embedding model $M_g(G) : G \to z(G)$, and (3) the final prediction model $M_o(h(\vec{x}), z(G)) : [h(\vec{x}); z(G)] \to \{P, N, O\}$. Hence, $h(\vec{x}) = h^{(L)}(\vec{x})$ for the TNET model and $h((\vec{x}) = h_{[CLS]}(\vec{x})$ for the BERT model. $z(G) = z(u)$ (section 3.2) and $M_o(h(\vec{x}), z(G))$ is a feedforward neural network with softmax output.

For generating global explanations from text features $h((\vec{x})$ we build on the method proposed in [40]. Let $\theta_{con}$ be the matrix of $k$-concept vectors used for explaining salient features of prediction, and $\theta_{dec}$ be the parameters of the decoder network $g$, which reconstructs the original vectors from concept embeddings $\theta_{con}^T h(\vec{x})$. The parameters of the text explanation model $E_t$ ($\theta_{con}, \theta_{dec}$) are learned by minimizing the regularized negative log-likelihood function:

$$\mathcal{L}_{E_t}(\theta_{con}, \theta_{dec}) = \mathbb{E}_{\vec{x}}[-\log P(M_o(g(\theta_{con}^T h(\vec{x}), \theta_{dec}), \; z(G))] + \mathbf{R}(\theta_{con})$$

where, $R(\theta_{con})$ is the diversity regularity between concepts described in [40]. Explanation words are extracted by selecting top scoring words $x_i$ according to the score $\max_k(\theta_{con}^T h(\vec{x_i})$. For the graph explanation model $E_g$ with parameters $\theta_g$, we use parameterized graph explainer (PGExplainer) model [18], minimizing the entropy loss:

$$\mathcal{L}_{E_g}(\theta_g) = \mathbb{E}_{G_S \sim q(\theta_g)}[H(M_o(h(\vec{x}), M_g(G)) \mid M_o(h(\vec{x}), M_g(G_S)))]$$

where $G_S$ is explanation subgraph sampled from a distribution $q$ parameterized by $\theta_g$.

A key challenge in generating multi-modal explanations is to identify the significance of individual modes. To this end, we define the significance variables for graph and text mode $s_g, s_t \in [0, 1]$, respectively. Given the explanations provided by unimodal explanation models, we generate the perturbed input text $\tilde{x}$ by removing the explanation words from $\vec{x}$, and the perturbed input graph $\tilde{G}$ by removing vertices and induced edges from the explanation subgraph. The significance labels are set as: $s_t = 1$ if

$M_o(h(\tilde{x}), z(G)) \neq M_o(h(\vec{x}), z(G))$ and 0 otherwise, and analogously for $s_g$. We also train feed-forward neural networks with sigmoid activation for predicting $s_t, s_g$ from inputs $\vec{x}, z(G)$: $s_t = S_t(\vec{x}, z(G))$, and $s_g = S_g(\vec{x}, z(G))$. Hence the joint multi-modal explanation loss is given as:

$$\mathcal{L}_{mm}(\theta_{con}, \theta_{dec}\theta_g, S_t, S_g) = S_t * \mathcal{L}_{E_t} + S_g * \mathcal{L}_{E_g} + \lambda(L(s_t, S_t) + L((s_g, S_g))$$

The first two terms are significance weighted explanation losses for individual modes, and the last two terms are binary cross-entropy losses for matching significance predictors to respective significance values generated by the unimodal explanation predictors. $\lambda$ controls the weightage given to initial unimodal significance values.

## 4 EXPERIMENTS

In this section, we report experimental results to empirically ascertain whether the proposed models indeed perform better than the existing state of the art methods.

### 4.1 Experimental Setup

**Datasets and baselines**: We evaluate our proposed models on the three benchmark datasets: LAPTOP and REST datasets from SemEval 2014 Task 4 subtask 2 [24] which contains reviews from Laptop and Restaurant domain respectively and the TWITTER dataset [7] containing Twitter posts. For TNet-based models, we perform the same prepossessing procedure as done in [15]. We compare results of our proposed models with state-of-the-art methods reported in table 2.

**Aspect disambiguation and KG Embedding**: For each aspect in the dataset $\mathcal{D}$ mentioned above, we disambiguate its corresponding entity in the knowledge graph using the *wikifier* API [2]. We use hierarchical Louvain graph clustering [1] algorithm for clustering the KG and constructing the weighted cluster graph $G_C(V_C, E_C, W_C)$(ref section 3.2.1). Statistics of the knowledge graph and its corresponding cluster graph and sub-graphs are shown in Table 1. For training entity sub-graph and weighted cluster graph embedding, we use GraphSAGE mean as aggregate function. For training GraphSage [9], we sample 25 nodes for layer 1 and 10 nodes for layer 2 using a random walk. The output hidden representation dimension is set as 50, and the number of negative samples $Q$ taken as 5. Default values are used for all other parameters.

**ALSC and probing function training**: For TNet-based models, we used 20% randomly held-out training data as the development set. We train the model for 100 epochs and select the model corresponding to the maximum development-set accuracy. Following [15], we use the same set of hyperparameters and report the maximum accuracy obtained on the test set over multiple runs. For training of BERT-based models, we use the procedure suggested in Rietzler et al. [27], for both pre-training and fine-tuning. For end-to-end training of ALSC with entity embedding generation, we use Adam optimizer with a learning rate of $3 \cdot 10^{-5}$, batch size of 512 for GraphSAGE-based entity embedding generation and 32 for ALSC task, number of epochs as 7. All the other hyper-parameters in GraphSAGE based entity embedding generation and ALSC task follow the same values in individual training. For training of the probing function $B$, we use Adam optimizer with a learning rate of $1 \cdot 10^{-5}$, the batch size of 128, the number of epochs as 100,

the probe dimension $dim_B$ as 100, and the regularization rate as 0.01. We use the same evaluation procedure suggested in Rietzler et al. [27], e.g. we conducted 9 runs with different initializations for all the experiments, and reported the average performance on the given test set. The model with the best training error over 7 epochs is taken as the final model for all runs.

**Table 1: Statistics of knowledge graph, weighted cluster graph and entity relation sub-graphs.**

| Knowledge Graph Embedding | | | |
|---|---|---|---|
| #Edges | #Nodes | #Clusters | Max. inter -cluster degree |
| 173068197 | 22504204 | 606 | 341 |
| Sub-graph Embedding | | | |
| Dataset | #Nodes | #Edges | Max. Node Degree |
| LAPTOP | 785 | 4477 | 107 |
| REST | 1031 | 7305 | 136 |
| TWITTER | 120 | 429 | 40 |

## 4.2 Comparison of ALSC Models

Table 2 reports all baseline and proposed models' performance, using two standard metrics: macro-accuracy (**ACC**) and macro-F1 score (**Macro-F1**). We observe that AR-BERT-idd outperforms all the other models on the REST and TWITTER dataset, and AR-BERT-S-idd outperforms all the other models for the LAPTOP dataset, in terms of both Accuracy and Macro-Averaged F1 scores. Hence, we can conclude that *representation of relations between aspect entities helps in training better ALSC models*. The improvement in the performance by the proposed AR-BERT-idd and AR-BERT-S-idd models over other BERT-based baseline models imply that the DBpedia knowledge graph encodes information which supplements the information contained in BERT embeddings of the aspect terms.We also note that BERT based baseline models, e.g. BERT-ADA, perform better than other models, e.g. TNET, as they utilize the context-sensitive word embeddings fine-tuned on domain-related datasets.
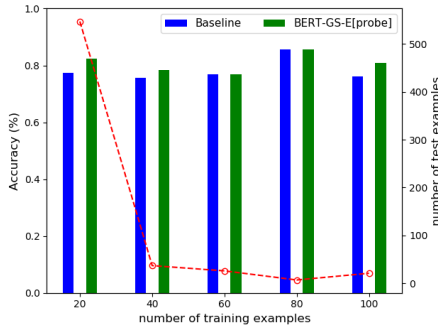


**Figure 2: Effect of number of training datapoints**

**Effect of training data scarcity**: Figure 2 reports the accuracy of the baseline model (BERT-ADA, blue bar) and the proposed model (green bar), for all test aspects in the LAPTOP dataset. The test

aspects are bucketed according to their training data counts, and the bars report average accuracy for all aspects in the buckets. We can see that for aspects that have $0 - 20$ training points, the proposed method outperforms the baseline. Hence, we conclude that for aspects with a low number of training data points, the proposed method improves the performance of ALSC by borrowing information from nearby aspects in the KG. The red line shows the number of test data points for each of the buckets. We find that a large fraction of test aspects have fewer than 20 training data points.

**Error Analysis**: Table 3 shows the confusion matrices of predictions of AR-TNet-idd and AR-BERT[idd] w.r.t. their respective baseline models on the three datasets. The top-left and bottom-right values report the number of correctly classified or misclassified examples by both methods in each sub-matrix. We can see that the proposed models do not induce any new errors which were not present in the respective baselines. Finally, we see that the bottom left entries in each table that report the new corrects (examples classified wrongly by existing methods but are classified correctly by the proposed methods) are much higher. Thus, we conclude that the new technique is an improvement over the old methods.

**Anecdotal examples**: Table 5 illustrates a few examples misclassified by BERT-ADA and correctly predicted by AR-BERT-idd. Aspects in the sentences are marked in the bracket with corresponding sentiment labels in subscript. Context words captured by the models for corresponding predictions are extracted using the context explainer of multi-modal explanation generation model, and the most important aspect entity (if any) is extracted using graph explainer of multi-modal explanation model (for instances where both context and graph information is required to predict sentiment polarity correctly, the most important aspect entity is the node in the explanation sub-graph with highest training example). Context explanations verify that BERT-ADA captures context words that are always semantically related to the aspects e.g. 'high' w.r.t. aspect 'price', whereas AR-BERT-idd considers aspect as an entity and tries to captures context words associated with that entity or the most important aspect entity from the given context. For each dataset, the distribution of test examples in each mode of importance of multi-modal explanation is given in Table 6.

## 4.3 Incorrect Disambiguation Detection

In this section, we demonstrate the effectiveness of our probing function for incorrect disambiguation detection. We categorized the aspects into 3 categories based on the disambiguation by wikifier: (1) **unknown (unk)** where there was no entity found, (2) **correct disambiguation (cd)** where the disambiguated aspect was mapped to the correct entity, and (3) **incorrect disambiguation (id)** where the disambiguated aspect was mapped to an incorrect entity, based on manual annotation. Table 4 shows the number of incorrectly classified **(ic)** examples (by the ALSC model) in each disambiguation category out of the total number of examples in that category ($\#ic/\#Total$). We see that compared to AR-BERT (wo end-to-end) and AR-BERT, AR-BERT-idd has significantly fewer incorrectly classified for the **unknown** and **incorrect disambiguation** categories. For the **correct disambiguation** category, all methods

**Table 2: Experiment results on various datasets(%). The marker * refers to p-value <0.01 when comparing with respective baselines. % in bracket of best performing models implies overall gain wrt. its' baselines.**

| Model | LAPTOP | | REST | | TWITTER | |
|---|---|---|---|---|---|---|
| | ACC | Macro-F1 | ACC | Macro-F1 | ACC | Macro-F1 |
| **Baseline models for ALSC** | | | | | | |
| TNet [15] | 76.33 | 71.27 | 79.64 | 70.20 | 78.17 | 77.17 |
| BERT-base [6] | 77.69 | 72.60 | 84.92 | 76.93 | 78.81 | 77.94 |
| SDGCN-BERT [45] | 81.35 | 78.34 | 83.57 | 76.47 | 78.54 | 77.72 |
| BERT-ADA [27] | 80.25 | 75.77 | 87.89 | 81.05 | 78.90 | 77.97 |
| **ALSC with Aspect relation incorporation** | | | | | | |
| AR-TNet | 78.80$^\star$ | 73.87$^\star$ | 83.40$^\star$ | 73.91$^\star$ | 80.52$^\star$ | 79.79$^\star$ |
| AR-BERT | 81.73$^\star$ | 77.07$^\star$ | 89.38$^\star$ | 82.47$^\star$ | 80.91$^\star$ | 80.15$^\star$ |
| AR-BERT-S | 82.37$^\star$ | 79.21$^\star$ | 85.27$^\star$ | 78.07$^\star$ | 79.67$^\star$ | 78.89$^\star$ |
| **ALSC with Aspect relation and incorrect disambiguation detection** | | | | | | |
| AR-TNet-idd | 80.09$^\star$ | 75.11$^\star$ | 84.64$^\star$ | 75.17$^\star$ | 81.64$^\star$ | 80.84$^\star$ |
| AR-BERT-idd | 82.91$^\star$ | 78.31$^\star$ | **90.62$^\star$** | **83.81$^\star$** | **82.08$^\star$** | **81.21$^\star$** |
| | | | (+3.11%) | (+3.40%) | (+4.03%) | (+4.15%) |
| AR-BERT-S-idd | **83.62$^\star$** | **80.43$^\star$** | 86.61$^\star$ | 79.37$^\star$ | 80.86$^\star$ | 80.03$^\star$ |
| | (+2.79%) | (+2.67%) | | | | |
| **Results with explanation concepts** | | | | | | |
| AR-BERT-idd ($h(x) = g(\theta_{con}^T h(\vec{x}), \theta_{dec})$) | - | - | 90.18 | 83.39 | 81.50 | 80.64 |
| AR-BERT-S-idd ($h(x) = g(\theta_{con}^T h(\vec{x}), \theta_{dec})$) | 83.38 | 80.25 | - | - | - | - |
| AR-BERT-idd ($G = G_S$) | - | - | 90.36 | 83.57 | 81.94 | 81.07 |
| AR-BERT-S-idd ($G = G_S$) | 83.54 | 80.41 | - | - | - | - |
| **Results with ablations of explanations** | | | | | | |
| AR-BERT-idd (with $\tilde{x}$) | - | - | 34.13 | 33.25 | 33.81 | 33.17 |
| AR-BERT-S-idd (with $\tilde{x}$) | 34.48 | 32.43 | - | - | - | - |
| AR-BERT-idd (with $\tilde{G}$) | - | - | 88.01 | 81.11 | 78.75 | 77.67 |
| AR-BERT-S-idd (with $\tilde{G}$) | 81.09 | 78.15 | - | - | - | - |

**Table 3: Confusion matrices of predictions of AR-TNet-idd vs TNet and AR-BERT-idd vs BERT-ADA w.r.t. correct and incorrect classification**

| Baseline | AR-TNet [idd] | | AR-BERT [idd] | |
|---|---|---|---|---|
| Prediction | Correct | Incorrect | Correct | Incorrect |
| **LAPTOP** | | | | |
| Correct | 487 | 0 | 512 | 0 |
| Incorrect | 24 | 127 | 17 | 109 |
| **REST** | | | | |
| Correct | 892 | 0 | 984 | 0 |
| Incorrect | 56 | 172 | 31 | 105 |
| **TWITTER** | | | | |
| Correct | 541 | 0 | 547 | 0 |
| Incorrect | 24 | 127 | 22 | 124 |

**Table 4: Fraction of incorrectly predicted examples in disambiguation categories.**

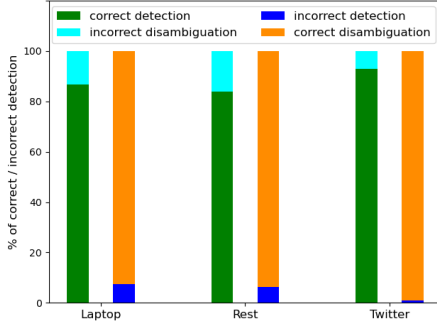| | # i_c / # u_d | # i_c / # i_d | # i_c / # c_d |
|---|---|---|---|
| **AR-BERT** | | | |
| LAPTOP | 8 / 10 | 47 / 53 | 62 / 575 |
| REST | 12 / 14 | 76 / 87 | 31 / 1019 |
| TWITTER | 2 / 2 | 8 / 14 | 122 / 676 |
| **AR-BERT-idd** | | | |
| LAPTOP | 6 / 10 | 41 / 53 | 62 / 575 |
| REST | 10 / 14 | 65 / 87 | 30 / 1019 |
| TWITTER | 0 / 2 | 0 / 14 | 124 / 676 |

our incorrect disambiguation detection method shows excellent performance, while also being highly scalable.

## 4.4 Effectiveness of Multi-modal Explanations

In this section, we report experimental results studying the effectiveness of multi-modal explanations predicted by the proposed model. Table 2, section *Results with explanation concepts* reports the predictive performance of ALSC models with explanations as input. The top two rows use the text concept embeddings as input, while the next two rows use the explanation subgraph as the input. We note that in all these cases the prediction performance remains similar to the original model, with the maximum accuracy difference

have the similar fraction of misclassification, which is much lower than the other two categories.

Figure 3 reports the percentage **i_d** which were detected correctly (left bar), and the **c_d** which were marked incorrectly (right bar) by the probing scheme. It can be seen that more than 80% of **i_d** examples have been correctly detected, and less than 5% of **c_d** examples have been wrongly flagged. Hence, we conclude that

**Table 5: Examples of mistakes by BERT-ADA which were correctly predicted by AR-BERT-idd. Red and green backgrounds indicate context explanations from multi-modal explanation model for BERT-ADA and AR-BERT-idd.**

| Sentence | BERT -ADA | AR-BERT -idd | Graph-mode Explanation |
|---|---|---|---|
| **LAPTOP** | | | |
| However, I can refute that [OSX]$_{NEG}$ is "FAST". | POS | NEG | Hard_disk_drive (39) |
| From the speed to the multi touch gestures this operating system beats [Windows]$_{NEG}$ easily. | POS | NEG | Software (64) |
| Did not enjoy the new [Windows 8]$_{NEG}$ and touchscreen functions. | NEG | NEG | - |
| **REST** | | | |
| Anywhere else, the [prices]$_{POS}$ would be 3x as high! | NEG | POS | Service (140) |
| A beautiful atmosphere, perfect for [drinks]$_{NEU}$ and/or appetizers. | POS | NEU | Food (365) |
| The [bread]$_{POS}$ is top notch as well. | POS | POS | - |
| **TWITTER** | | | |
| noobus Turns out [Snoop Dogg]$_{NEG}$ is actually pretty funny . | POS | NEG | Barack_Obama (343) |
| just got hold of an [Ipod]$_{NEU}$ . . it will be fun learning how to use it on the bus trip to canberra this monday | POS | NEU | IPhone (168) |
| Is it just me , or does [John Boehner]$_{NEG}$ sound like a newsman ? Sounds like he belongs on CBS Nightly News . | NEU | NEU | - |



**Figure 3: Percentage of correct / incorrect detection of disambiguation**

**Table 6: Distribution of test examples in each mode of multi-modal explanation with mode importance**

| Important Mode | LAPTOP | REST | TWITTER |
|---|---|---|---|
| **Text** ($S_t = 1$ and $S_g = 0$) | 300 (0) | 588 (0) | 355 (0) |
| **Graph** ($S_t = 0$ and $S_g = 1$) | 34 (0) | 45 (0) | 49 (0) |
| **Both** ($S_t = 1$ and $S_g = 1$) | 234 (15) | 427 (31) | 216 (22) |
| **None** ($S_t = 0$ and $S_g = 0$) | 70 (0) | 60 (0) | 72 (0) |
| **Total** | 638 | 1120 | 692 |

of $< 0.5\%$ and maximum macro-F1 score difference of $< 1.5\%$. This demonstrates the effectiveness of the explanation models. We also notice that predicting using the text concept embeddings results in marginally lower performance compared to predicting using graph explanation, which is expected due to the overall higher importance of the text mode. Table 2, section-*Results with ablations with explanations* further demonstates the effectiveness of explanations by reporting the performance of the ALSC model with perturbations $\tilde{x}$ and $\tilde{G}$ removing the explanations of text and graph modes as

input. The results (top two rows) clearly demonstrate that the text mode explanation include an important portion of the text, whose removal causes drastic drop in model performance.

Table 6 analyzes the effectiveness of prediction of important modes of explanation. Text mode implies $S_t = 1$ and $S_g = 0$, and other modes are defined analogously. We note that "Text" and "Both" mode cover a majority of the test examples, which is expected. Graph mode covers a miniscule number of examples for which the ALSC prediction can be made by simply observing the aspect graph nodes. Finally, we observe a small number of examples for which the explanation model predicts that none of the modes are important. A majority of these examples are ones where the ALSC model erroneously predicts the ALSC labels, due to which there is no possibility of improvement using any of the modes. The numbers in brackets in table 6 show the number of cases where there has been an improvement in label prediction by the ALSC model due to incorporation of aspect relations from graphs. Note that all the improvements have been in the cases of examples for which the explanation model predicts that both modes are important. This further demonstrates the effectiveness of the mode prediction models. Table 5 reports anecdotal examples from the three datasets along with text mode and graph mode explanations provided by the explanation model. Examples for each dataset shows examples with improvement in ALSC prediction due to incorporation of graphs. The entities corresponding to aspect graph nodes, and count of neighboring nodes in the aspect graph are shown in the graph mode explanation, which are intuitive.

## 5 CONCLUSIONS

In this paper, we present a scalable technique for incorporating aspect relations from large knowledge graphs, into state of the art deep learning based ALSC models. The resulting algorithm - AR-BERT, along with a novel incorrect disambiguation detection technique, results in consistent and significant improvements in

ALSC performance on all benchmark datasets. This work also reports the first algorithm for multi-modal explanation generation across textual and graph data.

# REFERENCES

[1] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. 2008. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* 2008, 10 (Oct 2008), P10008. `https://doi.org/10.1088/1742-5468/2008/10/p10008`

[2] Janez Brank, Gregor Leban, and Marko Grobelnik. 2018. Semantic Annotation of Documents Based on Wikipedia Concepts. *Informatica (Slovenia)* 42, 1 (2018). `http://www.informatica.si/index.php/informatica/article/view/2228`

[3] Samuel Broscheit. 2019. Investigating Entity Knowledge in BERT with Simple Neural End-To-End Entity Linking. In *Proceedings of the 23rd Conference on Computational Natural Language Learning, CoNLL 2019, Hong Kong, China, November 3-4, 2019*, Mohit Bansal and Aline Villavicencio (Eds.). Association for Computational Linguistics, 677–685. `https://doi.org/10.18653/v1/K19-1063`

[4] Peng Chen, Zhongqian Sun, Lidong Bing, and Wei Yang. 2017. Recurrent Attention Network on Memory for Aspect Sentiment Analysis. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*. 452–461. `https://www.aclweb.org/anthology/D17-1047/`

[5] Junqi Dai, Hang Yan, Tianxiang Sun, Pengfei Liu, and Xipeng Qiu. 2021. Does syntax matter? A strong baseline for Aspect-based Sentiment Analysis with RoBERTa. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, 1816–1829. `https://doi.org/10.18653/v1/2021.naacl-main.146`

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. `https://doi.org/10.18653/v1/n19-1423`

[7] Li Dong, Furu Wei, Chuanqi Tan, Duyu Tang, Ming Zhou, and Ke Xu. 2014. Adaptive Recursive Neural Network for Target-dependent Twitter Sentiment Classification. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 2: Short Papers*. 49–54. `https://www.aclweb.org/anthology/P14-2009/`

[8] Chunning Du, Haifeng Sun, Jingyu Wang, Qi Qi, Jianxin Liao, Tong Xu, and Ming Liu. 2019. Capsule Network with Interactive Attention for Aspect-Level Sentiment Classification. 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 2019, Hong Kong, China, November 3-7, 2019*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). Association for Computational Linguistics, 5488–5497. `https://doi.org/10.18653/v1/D19-1551`

[9] William L. Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*. 1024–1034. `http://papers.nips.cc/paper/6703-inductive-representation-learning-on-large-graphs`

[10] John Hewitt and Christopher D. Manning. 2019. A Structural Probe for Finding Syntax in Word Representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4129–4138. `https://doi.org/10.18653/v1/n19-1419`

[11] Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Seattle, Washington, USA, August 22-25, 2004*, Won Kim, Ron Kohavi, Johannes Gehrke, and William DuMouchel (Eds.). ACM, 168–177. `https://doi.org/10.1145/1014052.1014073`

[12] Bin Jiang, Jing Hou, Wanyue Zhou, Chao Yang, Shihan Wang, and Liang Pang. 2020. METNet: A Mutual Enhanced Transformation Network for Aspect-based Sentiment Analysis. In *Proceedings of the 28th International Conference on Computational Linguistics*. 162–172.

[13] Atsushi Kanehira, Kentaro Takemoto, Sho Inayoshi, and Tatsuya Harada. 2019. Multimodal Explanations by Predicting Counterfactuality in Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[14] Rijula Kar, Susmija Reddy, Sourangshu Bhattacharya, Anirban Dasgupta, and Soumen Chakrabarti. 2018. Task-Specific Representation Learning for Web-Scale Entity Disambiguation. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, Sheila A. McIlraith and Kilian Q. Weinberger (Eds.). AAAI Press, 5812–5819. `https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17281`

[15] Xin Li, Lidong Bing, Wai Lam, and Bei Shi. 2018. Transformation Networks for Target-Oriented Sentiment Classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*. 946–956. `https://doi.org/10.18653/v1/P18-1087`

[16] Bin Liang, Rongdi Yin, Lin Gui, Jiachen Du, and Ruifeng Xu. 2020. Jointly Learning Aspect-Focused and Inter-Aspect Relations with Graph Convolutional Networks for Aspect Sentiment Analysis. In *Proceedings of the 28th International Conference on Computational Linguistics*. 150–161.

[17] Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. K-bert: Enabling language representation with knowledge graph. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 2901–2908.

[18] Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang. 2020. Parameterized explainer for graph neural network. *arXiv preprint arXiv:2011.04573* (2020).

[19] Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. 2017. 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*. 4068–4074. `https://doi.org/10.24963/ijcai.2017/568`

[20] Navonil Majumder, Soujanya Poria, Alexander Gelbukh, Md Shad Akhtar, Erik Cambria, and Asif Ekbal. 2018. IARM: Inter-aspect relation modeling with memory networks in aspect-based sentiment analysis. In *Proceedings of the 2018 conference on empirical methods in natural language processing*. 3402–3411.

[21] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8779–8788.

[22] Matthew E Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A Smith. 2019. Knowledge Enhanced Contextual Word Representations. 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)*. 43–54.

[23] Nina Poerner, Ulli Waltinger, and Hinrich Schütze. 2020. E-BERT: Efficient-Yet-Effective Entity Embeddings for BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*. 803–818.

[24] Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. SemEval-2014 Task 4: Aspect Based Sentiment Analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation, SemEval@COLING 2014, Dublin, Ireland, August 23-24, 2014*. 27–35. `https://www.aclweb.org/anthology/S14-2004/`

[25] Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B. Viégas, Andy Coenen, Adam Pearce, and Been Kim. 2019. Visualizing and Measuring the Geometry of BERT. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 8592–8600. `http://papers.nips.cc/paper/9065-visualizing-and-measuring-the-geometry-of-bert`

[26] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*. 1135–1144.

[27] Alexander Rietzler, Sebastian Stabinger, Paul Opitz, and Stefan Engl. 2020. Adapt or Get Left Behind: Domain Adaptation through BERT Language Model Fine-tuning for Aspect-Target Sentiment Classification. In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asunción Moreno, Jan Odijk, and Stelios Piperidis (Eds.). European Language Resources Association, 4933–4941. `https://www.aclweb.org/anthology/2020.lrec-1.607/`

[28] Emanuel H Silva and Ricardo M Marcacini. [n. d.]. Aspect-based Sentiment Analysis using BERT with Disentangled Attention. ([n. d.]).

[29] Chi Sun, Luyao Huang, and Xipeng Qiu. 2019. Utilizing BERT for Aspect-Based Sentiment Analysis via Constructing Auxiliary Sentence. In *Proceedings of the*

*2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 380–385. https://doi.org/10.18653/v1/n19-1035

[30] Tianxiang Sun, Yunfan Shao, Xipeng Qiu, Qipeng Guo, Yaru Hu, Xuan-Jing Huang, and Zheng Zhang. 2020. CoLAKE: Contextualized Language and Knowledge Embedding. In *Proceedings of the 28th International Conference on Computational Linguistics*. 3660–3670.

[31] Duyu Tang, Bing Qin, and Ting Liu. 2016. Aspect Level Sentiment Classification with Deep Memory Network. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*. 214–224. https://www.aclweb.org/anthology/D16-1021/

[32] Hao Tang, Donghong Ji, Chenliang Li, and Qiji Zhou. 2020. Dependency graph enhanced dual-transformer structure for aspect-based sentiment classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 6578–6588.

[33] Jialong Tang, Ziyao Lu, Jinsong Su, Yubin Ge, Linfeng Song, Le Sun, and Jiebo Luo. 2019. Progressive Self-Supervised Attention Learning for Aspect-Level Sentiment Analysis. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, Anna Korhonen, David R. Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, 557–566. https://doi.org/10.18653/v1/p19-1053

[34] Shuai Wang, Sahisnu Mazumder, Bing Liu, Mianwei Zhou, and Yi Chang. 2018. Target-Sensitive Memory Networks for Aspect Sentiment Classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*. 957–967. https://doi.org/10.18653/v1/P18-1088

[35] Yequan Wang, Minlie Huang, Xiaoyan Zhu, and Li Zhao. 2016. Attention-based LSTM for Aspect-level Sentiment Classification. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*. 606–615. https://www.aclweb.org/anthology/D16-1058/

[36] Bowen Xing and Ivor W. Tsang. 2021. Understand me, if you refer to Aspect Knowledge: Knowledge-aware Gated Recurrent Memory Network. *CoRR* abs/2108.02352 (2021). arXiv:2108.02352 https://arxiv.org/abs/2108.02352

[37] Lu Xu, Lidong Bing, Wei Lu, and Fei Huang. 2020. Aspect Sentiment Classification with Aspect-Specific Opinion Spans. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 3561–3567. https://doi.org/10.18653/v1/2020.emnlp-main.288

[38] Heng Yang, Biqing Zeng, Mayi Xu, and Tianxing Wang. 2021. Back to Reality: Leveraging Pattern-driven Modeling to Enable Affordable Sentiment Dependency Learning. *CoRR* abs/2110.08604 (2021). arXiv:2110.08604 https://arxiv.org/abs/2110.08604

[39] Heng Yang, Biqing Zeng, Jianhao Yang, Youwei Song, and Ruyang Xu. 2021. A multi-task learning model for Chinese-oriented aspect polarity classification and aspect term extraction. *Neurocomputing* 419 (2021), 344–356. https://doi.org/10.1016/j.neucom.2020.08.001

[40] Chih-Kuan Yeh, Been Kim, Sercan O Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. 2019. On completeness-aware concept-based explanations in deep neural networks. *arXiv preprint arXiv:1910.07969* (2019).

[41] Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. Gnnexplainer: Generating explanations for graph neural networks. *Advances in neural information processing systems* 32 (2019), 9240.

[42] Chen Zhang, Qiuchi Li, and Dawei Song. 2019. 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)*. 4560–4570.

[43] Yue Zhang and Jiangming Liu. 2017. Attention Modeling for Targeted Sentiment. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 2: Short Papers*. 572–577. https://www.aclweb.org/anthology/E17-2091/

[44] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 1441–1451.

[45] Pinlong Zhao, Linlin Hou, and Ou Wu. 2020. Modeling sentiment dependencies with graph convolutional networks for aspect-level sentiment classification. *Knowl. Based Syst.* 193 (2020), 105443. https://doi.org/10.1016/j.knosys.2019.105443

**Table 7: Statistics of datasets.**

|         | # Positive | # Negative | # Neutral |
|---------|-----------|-----------|-----------|
| LAPTOP  |           |           |           |
| Train   | 987       | 866       | 460       |
| Test    | 341       | 128       | 169       |
| REST    |           |           |           |
| Train   | 2164      | 805       | 633       |
| Test    | 728       | 196       | 196       |
| TWITTER |           |           |           |
| Train   | 1567      | 1563      | 3127      |
| Test    | 174       | 174       | 346       |

**Table 8: Performance comparison of our proposed models with other competitive models.**

| Model | LAPTOP | REST |
|-------|--------|------|
| LSA+DeBERTa-V3-Large | **86.21** | **91.07** |
| LCF-ATEPC | 82.29 | 90.18 |
| ABSA-DeBERTa | 82.76 | 89.46 |
| RoBERTa+MLP | 83.78 | 87.37 |
| KaGRMN-DSG | 81.87 | 87.35 |
| AR-BERT-idd | 82.91 | 90.62 |
| AR-BERT-S-idd | 83.62 | 86.61 |

## A APPENDIX

### A.1 Data Statistics

The statistics of the three datasets are given in Table 7.

### A.2 Implementation Details

Our implementation is based on Tensorflow 1.X and all our experiments were executed on Quadro P5000 Single core GPU (16278MiB) with CUDA version: 10.0.

*A.2.1 Aspect disambiguation.* We pass the sentence to the Wikifier using 'POST' method and extract the entity from the annotation of the maximum sub-string of the aspect. If no annotation is found we link the retrieved entity as 'Unknown'.

*A.2.2 Knowledge Graph Clustering.* We extract the maximum weakly connected component of the DBPedia Knowledge Graph (KG) by converting the KG into a SNAP undirected graph[2] (fraction of nodes in maximum weakly connected component: 0.99983). We then cluster the extracted maximum weakly connected component using hierarchical Louvain graph clustering algorithm. This algorithm returns a hierarchy tree with 5 levels in less than 40 minutes and the nodes in each level is as follows:

- level 0: 22504204 nodes
- level 1: 227550 nodes
- level 2: 2938 nodes
- level 3: 769 nodes
- level 4: 606 nodes

*A.2.3 Domain Specific BERT pre-training.* Similar to [27], we pre-train the BERT model on the two domain specific publicly available

datasets: Amazon electronics reviews[3] for LAPTOP domain and Yelp restaurants dataset[4] for REST domain. We adopt the similar pre-processing steps of [27] by removing reviews with less than 2 sentences to enable Next Sentence Prediction (NSP) task of BERT language model. We also removes the reviews from Amazon review data which appear in the LAPTOP dataset to eliminate training bias towards those reviews. After pre-processing, we get around 1 million reviews for LAPTOP domain and we sample around 10 million reviews from pre-processed Yelp review dataset for REST domain. We run the pre-training of BERT model for 30 epochs and 3 epochs for LAPTOP and REST domain respectively to train the language model with significant large amount of data (equal for both domains).

### A.3 Other Competitive ALSC Models

In this section we compare our proposed models performance with other competitive ALSC models:

- **LSA+DeBERTa-V3-Large** [38] This work introduces sentiment pattern based sentiment dependency learning framework to model the sentiment dependency between the adjacent aspects to learn the aspect polarity of an aspect without explicit sentiment context.
- **LCF-ATEPC** [39] This work focuses on multi-task learning based aspect term extraction and aspect sentiment polarity classification by introducing local context focusing mechanism.
- **ABSA-DeBERTa** [28] This work introduces DeBERTa model (Decoding-enhanced BERT with Disentangled Attention) in the ALSC task and achieves competitive performance.
- **RoBERTa+MLP** [5] This work examines the effectiveness of induced tree from the Pre-trained language model and utilizes induced tree from fine-tuned RoBERTa on ALSC task to acheieve competitive performance.
- **KaGRMN-DSG** [36] This work utilizes both local and global syntactic representation of aspects combined with knowledge representation of aspects to enhance the overall aspect representation.

We compare our proposed models performance with other competitive ALSC models in Table 8. Our models stands 2nd and 3rd for the REST and LAPTOP domain respectively in the recent leader board for ALSC.

### A.4 Additional Results

The extended version of Table 2 is given in Table 9 where we report the performance of our proposed ALSC models (without end-to-end training). Table 10 illustrates anecdotal examples of multi-modal aspect extraction on the three datasets for different mode importance.

---

**Table 9: Extended experiment results on various datasets(%). The marker * refers to p-value <0.01 when comparing with respective baselines. % in bracket of best performing models implies overall gain wrt. its' baselines.**

| Model | LAPTOP | | REST | | TWITTER | |
|---|---|---|---|---|---|---|
| | ACC | Macro-F1 | ACC | Macro-F1 | ACC | Macro-F1 |
| **Baseline models for ALSC** | | | | | | |
| TNet [15] | 76.33 | 71.27 | 79.64 | 70.20 | 78.17 | 77.17 |
| TNet-ATT [33] | 77.62 | 73.84 | 81.53 | 72.90 | 78.61 | 77.72 |
| BERT-base [6] | 77.69 | 72.60 | 84.92 | 76.93 | 78.81 | 77.94 |
| SDGCN-BERT [45] | 81.35 | 78.34 | 83.57 | 76.47 | 78.54 | 77.72 |
| BERT-ADA [27] | 80.25 | 75.77 | 87.89 | 81.05 | 78.90 | 77.97 |
| **ALSC with Aspect relation incorporation** | | | | | | |
| AR-TNet (wo end-to-end) | 77.89★ | 72.96★ | 82.31★ | 72.97★ | 79.68★ | 78.83★ |
| AR-TNet | 78.80★ | 73.87★ | 83.40★ | 73.91★ | 80.52★ | 79.79★ |
| AR-BERT (wo end-to-end) | 80.87★ | 76.13★ | 88.21★ | 81.45★ | 79.83★ | 79.02★ |
| AR-BERT | 81.73★ | 77.07★ | 89.38★ | 82.47★ | 80.91★ | 80.15★ |
| AR-BERT-S (wo end-to-end) | 81.82★ | 78.75★ | 84.64★ | 77.34★ | 79.06★ | 78.36★ |
| AR-BERT-S | 82.37★ | 79.21★ | 85.27★ | 78.07★ | 79.67★ | 78.89★ |
| **ALSC with Aspect relation and incorrect disambiguation detection** | | | | | | |
| AR-TNet-idd | 80.09★ | 75.11★ | 84.64★ | 75.17★ | 81.64★ | 80.84★ |
| AR-BERT-idd | 82.91★ | 78.31★ | **90.62★** | **83.81★** | **82.08★** | **81.21★** |
| | | | (+3.11%) | (+3.40%) | (+4.03%) | (+4.15%) |
| AR-BERT-S-idd | **83.62★** | **80.43★** | 86.61★ | 79.37★ | 80.86★ | 80.03★ |
| | (+2.79%) | (+2.67%) | | | | |
| **Results with explanation concepts** | | | | | | |
| AR-BERT-idd ($h(x) = g(\theta_{con}^T h(\vec{x}), \theta_{dec})$) | - | - | 90.18 | 83.39 | 81.50 | 80.64 |
| AR-BERT-S-idd ($h(x) = g(\theta_{con}^T h(\vec{x}), \theta_{dec})$) | 83.38 | 80.25 | - | - | - | - |
| AR-BERT-idd ($G = G_S$) | - | - | 90.36 | 83.57 | 81.94 | 81.07 |
| AR-BERT-S-idd ($G = G_S$) | 83.54 | 80.41 | - | - | - | - |
| **Results with ablations of explanations** | | | | | | |
| AR-BERT-idd (with $\tilde{x}$) | - | - | 34.13 | 33.25 | 33.81 | 33.17 |
| AR-BERT-S-idd (with $\tilde{x}$) | 34.48 | 32.43 | - | - | - | - |
| AR-BERT-idd (with $\tilde{G}$) | - | - | 88.01 | 81.11 | 78.75 | 77.67 |
| AR-BERT-S-idd (with $\tilde{G}$) | 81.09 | 78.15 | - | - | - | - |

**Table 10: Examples of Multi-modal Explanation Extraction. Aspect in the parenthesis, actual label in the subscript, predicted label around the aspect with proper color code: green for positive, red for negative and yellow for neutral, context explanation tokens are in blue color.**

| Sentence | Aspect-Entity | Explanation Node |
|---|---|---|
| **Mode Importance: Context** | | |
| [**Boot time**]POS is super fast, around anywhere from 35 seconds to 1 minute. | - | - |
| The [**bread**]POS is top notch as well. | - | - |
| i like [**Britney Spears**]POS new song ... i wan na hear it now = -LRB- | - | - |
| **Mode Importance: Both** | | |
| Did not enjoy the new Windows 8 and [**touchscreen functions**]NEG. | **Touchscreen** | **Software(64)** |
| Did I mention that the [**coffee**]POS is OUTSTANDING? | **Coffee** | **Drink (35)** |
| RT jaimemorelli : I would love to see a nuanced comparison of [**Google**]NEU television vs . hooking up my television to a Mac Mini and buying a wireless keyboard | **Google** | **iPhone (168)** |
| **Mode Importance: Graph** | | |
| I would have given it 5 starts was it not for the fact that it had [**Windows 8**]NEG | **Windows_8** | **Software(64)** |
| The [**food**]NEU did take a few extra minutes to come, but the cute waiters' jokes and friendliness made up for it. | **Food** | **Food (365)** |
| [**Shaquille O'Neal**]NEG to miss 3rd straight playoff game \| The ... : shaquille o'neal will miss his third straight play | **Shaquille_O'Neal** | **Los_Angeles_Lakers (157)** |